

## N O T I C E

THIS DOCUMENT HAS BEEN REPRODUCED FROM  
MICROFICHE. ALTHOUGH IT IS RECOGNIZED THAT  
CERTAIN PORTIONS ARE ILLEGIBLE, IT IS BEING RELEASED  
IN THE INTEREST OF MAKING AVAILABLE AS MUCH  
INFORMATION AS POSSIBLE

"Made available under NASA sponsorship  
in the interest of early and wide dis-  
semination of Earth Resources Survey  
Program information and without liability  
for any use made thereof."

80-10315  
CR-163403  
SDSU-RSI-79-03

IMPACT OF CELL SIZE ON INVENTORY  
AND MAPPING ERRORS IN A CELLULAR  
GEOGRAPHIC INFORMATION SYSTEM

(E80-10315) IMPACT OF CELL SIZE ON  
INVENTORY AND MAPPING ERRORS IN A CELLULAR  
GEOGRAPHIC INFORMATION SYSTEM (South Dakota  
State Univ.) 96 p HC A05/MF A01 CSCL 08B

N80-33831

Unclas  
G3/43 00315

To

National Aeronautics Space Administration  
Washington, D.C.  
Grant Number NGL-42-003-007

Remote Sensing Institute  
South Dakota State University  
Brookings, South Dakota 57007

April 1979

Interim Technical Report

IMPACT OF CELL SIZE ON INVENTORY  
AND MAPPING ERRORS IN A CELLULAR  
GEOGRAPHIC INFORMATION SYSTEM

by

Michael E. Wehde

To

National Aeronautics Space Administration  
Washington, D.C.  
Grant Number NGL-42-003-007

Remote Sensing Institute  
South Dakota State University  
Brookings, South Dakota 57007

April 1979

## ABSTRACT

Cellularization of spatially distributed data for manipulation by computerized information systems gives rise to mapping and inventory errors in the data sets. The impact of cell size on these errors was investigated. Mapping error in particular was studied with changing cell size for individual mapping units and a four-square mile segment of soil survey data. The effect of grid position was analyzed and found insignificant for maps but highly significant for isolated mapping units. A modelable relationship between mapping error and cell size was observed for the map segment analyzed. Map data structure was also analyzed with an interboundary-distance distribution approach. Map data structure and the impact of cell size on that structure were observed. The existence of a model allowing prediction of mapping error based on map structure was hypothesized and two generations of models were tested under simplifying assumptions. Results of the modeling efforts are reviewed and application significance discussed.

## ACKNOWLEDGEMENTS

The author wishes to express appreciation to Dr. Victor I. Myers, Director of the Remote Sensing Institute and the Institute staff for enabling the research to be performed. The work was supported by the NASA research grant number NGL-42-003-007. The comments of Dr. Robert Finch of the Electrical Engineering Department were most helpful.

Consultations with Mr. Delmar Johnson, the Systems Programmer, at the SDSU computer facility were invaluable to the successful implementation of AREAS. The author is grateful to the following people: Lynette Nelson, for patience in typing repeated adjustments to the text; Mary Buckmiller and Kathy Kitzmiller, for graphic arts assistance in figure production; and Jan Griesenbrock, for timely production of photographic products during the course of the investigation and during the production of figures.

## TABLE OF CONTENTS

	PAGE
ABSTRACT . . . . .	i
ACKNOWLEDGEMENTS . . . . .	ii
TABLE OF CONTENTS . . . . .	iii
LIST OF FIGURES . . . . .	v
LIST OF TABLES . . . . .	vii
INTRODUCTION . . . . .	1
GEOGRAPHICAL INFORMATION SYSTEMS . . . . .	2
DATA BASE CHARACTERISTICS . . . . .	3
INFORMATION SYSTEM EXAMPLES . . . . .	6
SYSTEM STUDIES AND RECOMMENDATIONS . . . . .	8
SYSTEM COMPARISONS AND EVALUATIONS . . . . .	10
DATA ORGANIZATION . . . . .	10
PERFORMANCE EVALUATIONS . . . . .	14
GRID CELL SIZE SELECTION BASIS . . . . .	18
RESEARCH OBJECTIVES . . . . .	19
SYSTEM PERFORMANCE EXPERIMENTS . . . . .	20
THE DATA BASE . . . . .	21
PERFORMANCE EXPERIMENTS . . . . .	24
EXPERIMENTAL RESULTS . . . . .	25
PERFORMANCE RELATIONSHIPS . . . . .	30
AN ORIENTATION STUDY . . . . .	34
DATA CHARACTERIZATION ANALYSIS . . . . .	42
DISTRIBUTION OF SPANS . . . . .	42
SPANS VS CELL SIZE . . . . .	52
DATA CHARACTERIZATION . . . . .	54

## TABLE OF CONTENTS (CONT'D)

	PAGE
MAPPING ERROR AND THE SPAN DISTRIBUTION	
- A POSITIONAL AVERAGE MODEL . . . . .	55
THE SIZE AND ORIENTATION OF CELLS . . . . .	55
THE SPAN DISTRIBUTION AND MAPPING ERROR . . . . .	56
A POSITIONAL AVERAGE MODEL . . . . .	58
PREDICTION OF MAPPING ERROR . . . . .	62
MAPPING ERROR AND THE SPAN DISTRIBUTION	
- CORRECTION FOR SPAN ADJACENCIES . . . . .	63
SPAN ADJACENCY INFLUENCE . . . . .	64
DESCRIBING ADJACENT SPAN CORRECTIONS . . . . .	65
THE CORRECTION MATRIX . . . . .	66
IMPLEMENTING THE CORRECTION . . . . .	75
APPLYING THE CORRECTION . . . . .	78
SUMMARY, CONCLUSIONS AND RECOMMENDATIONS . . . . .	81
SUMMARY . . . . .	81
CONCLUSIONS . . . . .	81
RECOMMENDATIONS . . . . .	83
BIBLIOGRAPHY . . . . .	86

## LIST OF FIGURES

FIGURE		PAGE
1	The central-datum cell coding concept as related to mapping error . . . . .	16
2	The original map of the intensive-study data set . . . . .	22
3	The computer map of the "true" data set at 0.007 ha (0.0174 acre) cellular grid . . . . .	23
4	Performance evaluation processing diagram for an increased cell size . . . . .	25
5	The twelve maps compared for mapping accuracy versus cell size . . . . .	26
6	Representations of mapping error . . . . .	27
7	Mapping and inventory errors versus resolution number . . . . .	29
8	Inventory errors versus resolution number for selected mapping units . . . . .	31
9	Mapping errors versus resolution number for selected mapping units . . . . .	32
10	Spatial representations of the mapping units referred to in Figure 8 . . . . .	33
11	The effect of grid position on mapping error for a single circular mapping unit . . . . .	35
12	Mapping error versus resolution number for four single closed regions . . . . .	37
13	Growth of area error as cell size increases for a fixed interboundary map distance . . . . .	44
14	The span distribution reflects region size, orientation and regularity of shape . . . . .	46
15	Span distributions for the reference data set of resolution number one . . . . .	47
16	Span distributions for combined horizontal and vertical scans of the twelve different cell-sized data sets . . . . .	49



## LIST OF FIGURES (Cont'd)

FIGURE	PAGE
17 Span means versus resolution numbers . . . . .	52
18 Mean span distance versus resolution number . . . . .	54
19 Components of the proposed mathematical relationship between span distribution and mapping error vectors . . . . .	58
20 The procedure for observing $\bar{g}(n,m)$ by systematic sketches of the $m$ positions of a cell of size $m$ with respect to a span of size $n$ . . . . .	60
21 Predicted versus experimental mapping error . . . . .	62
22 Impact of an adjacent span on calculation of positional average mapping error . . . . .	64
23 One possible position of a six-unit cell with respect to a two-unit span . . . . .	66
24 One possible position of a six-unit cell with respect to a three-unit span . . . . .	67
25 One possible position of a seven-unit cell with respect to a three-unit span . . . . .	67
26 The experimentally observed mapping error compared to the positional-average model and the span-adjacency-corrected model . . . . .	80

## LIST OF TABLES

TABLE	PAGE
1 Mapping unit characteristics for mapping units referenced in Figures 3.7 and 3.8 . . . . .	34
2 The Kolmogorov-Smirnov test of common parent span distribution for the horizontal and vertical span distributions of the twelve data sets . . . . .	50
3 Positional average error fractions as a model of $\bar{g}(n,m)$ . . . . .	61
4 Correction expressions for selected m-n combinations . . . . .	70
5 Pattern of correction expressions of Table 4 . . . . .	71
6 Coefficients of $\bar{f}(n)$ for various m to generate first-order correction expressions for span adjacency effects . . . . .	72
7 Terms for even and odd m in the "k" entries of Table 5 . . . . .	73
8 The "T" entries of Table 5 . . . . .	74
9 Coefficients of k in the entries to the $\bar{c}(m,n)$ matrix . . . . .	76
10 Subtractive numbers for entries to the $\bar{c}(m,n)$ matrix . . . . .	77
11 Mapping error comparisons for the positional average model, the span adjacency correction, the corrected model and the experimental results . . . . .	79

## INTRODUCTION

Digital hardware evolution and software innovation have contributed to increased storage capacity, high speed processing and improved cost-performance ratio in data processing applications. Continuing improvements in storage, speed and costs of digital systems broaden the spectrum of applications. Present technological capabilities have well surpassed the minimum requirements for useful processing of geographically oriented, spatially distributed data sets.

Spatial data sets carry an added information dimension, i.e. location. Locational information allows interaction and display in a map reference framework. Methods for storage, retrieval and manipulation of such data sets, while preserving the location information, vary in complexity, accuracy and cost.

The inability of man to mentally cope with the complexity of processing or physically cope with the masses of data involved motivated investigations of computer processing. Other motivations have been noted. Nichols [1] pointed out that computerization of resource maps allow automated reduction of many levels of map information down to a few classes, suitabilities or potentials according to the specific purpose of the analysis and, thereby, restores human visualization to the problem. McDonald [2] cites an ever growing requirement for increased rates of delivery. Either situation alone would be cause for pursuit of computer capability.

Computer systems do have the capability. Over the last two decades software packages have proliferated in response to the combination of available hardware and specific problems to be handled. National and international examples which are identified later are LUNR, LUIS, MIADS, ORRMIS, GRIDS, CMS, NARIS, DIME, MLIS, GRDSR, MIDAS, GIST, FRIS, CGIS, MAP/MODEL, PIOS, and STORET [4]. Tomlinson [5] cites AUTOMAP, SYMAP, CAN-HYDRO, SYMVUU, GIMMS, OEM, NCC, ARDS, CENSUS, OBLIZ, and WWW. Bryant [6] offers IBIS as yet another alternate approach.

### Geographical Information Systems

The data base influences decisions about system development. As a significant factor in the design and operation of the information system, the data base must receive careful attention.

### Data Base Characteristics

All computer data sets are collections of data. The entries may be measurements of a variable, identifications of class membership or theme level, and at least implicitly the time or date of observation. The relevance of this time datum is dependent on the dynamics of the theme or variable.

Spatial or geographical information processing, however, deals in data sets, often termed data bases, which contain the theme or variate value, the time reference, and spatial location [7]. Steiner [7] lists the factors which determine the spatial characteristics of the data base as (A) the area span (geographical extent or coverage of the data base), (B) the spatial resolution (minimum size cell or spatial

unit represented), and (C) the spatial frequency (the occurrence of the resolution elements throughout the coverage, from a sampling system to a complete representation). He further points out that similar characteristics may be defined for the time or date datum, i.e. temporal span, temporal resolution, and temporal frequency. These characteristics of course only apply to data bases containing a theme or variable repeated at multiple time references as would be the case in studies of time dynamics.

Clearly the content variation in data bases discussed to this point is limited to multiple theme or variable data at a point in time, multiple time samples of a theme or variable, or both. The intended analyses guide formulation of the composition of the data base and also dictates to some degree the development of the information processing system. This is the first diverging influence which propagates new processing systems.

Determining data base content in light of purpose does not necessarily determine the measurement scale for the datum points. Steiner [7] describes several scales to choose from. First, a nominal scale indicates presence or absence of a theme. Second, an ordinal scale places observations in a size sequence, ranking or ordering of the theme or variable value. Third, a continuous scale is usually associated with a continuous variable and may be either interval or ratio type. Interval continuous scaling has an arbitrary zero and maintains data value differences, e.g. temperature measurements. Ratio continuous scaling has a physically meaningful absolute zero which then allows analysis by ratios of data values as well as differences

among values, e.g. rainfall measurement. In addition the conception of a thematic scale is an extension of nominal scaling to multiple levels not of an ordinal or continuous nature.

The application or intended analysis which guides selection of data base content variables or themes will also typically guide selection of the measurement scales. The variety of data scales, however, has greater impact on the processing system than the data base content discussed earlier. This is the second diverging influence which motivates creation of new processing systems.

In geographic information systems, data base design decisions of content and scale are subordinate to consideration of the alternatives for representation of the locational information. In fact geographic information systems are classified into types based primarily on the location reference element of the data base. Five of the most commonly recognized types according to approach to location reference are (1) the uniform grid which divides the space into an x-y cellular network, (2) the parcel in which spatial subdivisions arise in either natural or artificial (political) boundary context and are irregularly shaped, non-uniformly sized or both, (3) the area boundary in which addresses of the enclosing border are stored in conjunction with the data for the enclosed region, (4) the network in which lines connect nodes together in a spatial mesh or net, and (5) the point where spot spatial addressing places the data record into the overall spatial domain but the data values have no single uniform spatial coverage around the specific site [4].

Bryant [6] suggests a sixth type, image format, which is a sequential raster of pixels and greatly resembles the uniform grid in spatial philosophy but not necessarily in processing methodology.

Tomlinson [8] proposes four types of geographical data bases according to the handling of the locational identifiers. He suggests (1) external index, (2) coordinate reference, (3) arbitrary grid, and (4) explicit boundary. The first is identifiable by data sets containing no direct reference to location nor physical arrangement of the data which might imply spatial location. A separate index file contains the geographic locations and address pointers to identify the position of the data in the data base. The second system carries geographic coordinates together with the data entries. The third system divides the spatial domain into a regularly spaced, cellular network and the data are stored in such a manner as to preserve the relative location of these cells. The fourth system stores specific locational coordinates for the boundaries which enclose a homogeneous region. Although all of these types of data base accomplish the task of maintaining geographic or spatial reference, it is particularly important to note that only the latter two actually store and preserve boundary information. Hence only these two systems lend themselves to spatial displays of boundaries [8].

This preservation of boundary information for purposes of spatial display is precisely in line with a prime motivation of geographical information processing, i.e. mapping of information. For this reason only two types of information systems are commonly recognized. These are the grid encoding (cellular) and line encoding (polygon) approaches [9].

It is quite evident that the approach to handling the locational identifiers in a geographical data base has the greatest impact on the processing system of any of the factors discussed. This is the third and most strongly divergent influence on development of information systems.

With data base contents, analysis purposes, measurement scales, and geographical reference approaches varying from application to application and, furthermore, the idiosyncrasies of specific hardware configurations confounding the problem, the reason for the proliferation of data processing systems becomes apparent.

#### Information System Examples

Examples of systems employing the uniform grid approach to geographic reference are as follows: LUNR - Land Use and Natural Resources by the Office of Planning and Coordination for the State of New York and Cornell Center for Aerial Photographic Studies; LUIS - Land Use Information System by the University of Massachusetts; MIADS - Map Information Assembly and Display System by the U.S. Forest Service within the U.S. Department of Agriculture; ORRMIS - Oak Ridge Regional Modeling Information System by the Oak Ridge National Laboratory; GRIDS - Grid Relation Information Display System by the Southern California Regional Information Study; and CMS - Composite Mapping System by the Economic Development Administration within the U.S. Department of Commerce [4].

Examples of geographic information systems employing the parcel approach to location are as follows: NARIS - Natural



Resources Information System developed at Center for Advanced Computation at the University of Illinois for the Northeast Illinois Natural Resource Service Center; DIME - Dual Independent Map Encoding by the Bureau of Census within the U.S. Department of Commerce; MLIS - Minnesota Land Management Information System Study by the University of Minnesota; GRDSR - Geographically Referred Data Storage and Retrieval System by the Dominion Bureau of Statistics in Canada; GIST - Geographic Information System for the Office of the Mayor of New York City and FRIS by the Swedish Central Board for Real Estate Data [4].

Examples of geographic information systems employing the area boundary approach to spatial location are as follows: CGID - Canadian Geographic Information System by IBM corporation for the Canadian Department of Regional Economic Expansion; MAP/MODEL by the University of Oregon for the Bureau of Governmental Research and Service; and PIOS (no further information available). The lone example of network approach is STORET (no further information provided)[4].

IBIS - Image Based Information System is a newcomer in processing approach although the storage method is actually a fine cell or uniform grid basis [6].

The multitude of systems is clear evidence of the usefulness of computer technology in the area of geographic information processing. Furthermore, the many systems exist as a result of the influences discussed earlier which tend to generate new systems to meet new situations.

### System Studies and Recommendations

From a data processing or technological viewpoint a system is usually considered to be the combination and interaction of hardware and software. The typical consumer of geographic information would like to approach the system with only questions and analyses in mind, hoping to find both the processing capability and the requisite data sets available. Thus, from a user viewpoint the data base is a part of the information system.

Steiner [7] cites data base development as a five phase activity; data specification, data acquisition, data storage/retrieval/manipulation, data dissemination, and data applications. These activities, if included as a part of the information system, would promote continued development of new systems for new problems as the users and their data attempt to interact with the storage and processing capabilities of their hardware/software.

There are many opinions as to the content and developmental needs in generating useful geographic information systems.

Bryant and Zobrist [10] offer four criteria which they believe must be satisfied to make a geographic information system "useful". These criteria are (1) that point and area locations are provided, (2) that variable aggregation or sub-setting be possible, (3) that there be representation of spatial arrangement of the data, and (4) that data interface with mathematical and statistical analysis programs. These criteria are indicative of recognized needs for spatial manipulation and analytical capabilities in a useful system.

In a feasibility study for the Illinois Resources Information System, six areas were deemed important to the creation of a viable information system [4]. These were as follows: (1) that the system involve a large user community, (2) that point, network, and area data exist in one system, (3) that multiple data bases and resolutions be employed within the system, (4) that multiple problem-oriented user interfaces or terminal languages be provided, (5) that multiple data entry facilities spanning a variety of forms and formats of data be provided, and (6) that advanced graphics output capability be included.

The United States Geological Survey, in cooperation with the International Geographical Union's Commission on Geographical Data Sensing and Processing conducted a two-year study of spatial data systems. W.A. Radlinski, the Associate Director of USGS, in his opening address to the 1977 ASP-ACSM Annual Convention reported on specific findings of that study as follows [3]: (1) many seemingly simple tasks require development, (2) storage schemes need to consider large files up to terrabits ( $10^{12}$ ), (3) the digitization and edit function must be made more economical, (4) one system to handle point, line, area, network and image data is needed, (5) topological structure in data needs to be understood, and (6) storage/retrieval methods far exceed analytical/interpretive methods. His summary of cooperative developmental areas called for (1) establishing standards on scale, resolution, accuracy etc. by means of user survey, (2) sorting the potential or candidate data to be maintained in a system into priority by usefulness, (3) coordination of the

identification of theoretical research needs, (4) cooperative hardware development, (5) technology forecasting to influence system development planning, and (6) establishing a mechanism for institutional communication and coordination. His review and recommendations represent the broadest and most progressive noted in the literature.

Philips [9], although biased toward value of graphics, did make the significant observation that computer graphics must serve the role of reducing "impedance of the interface" between the data and the users.

Common threads of meaning in the recommendations cited appear to be (1) concern for the user and his interests and (2) concern for controlled, progressive, positive developmental effort in geographic information systems.

#### System Comparisons and Evaluations

Tomlinson [8] emphasized in his review of geographical data bases that only cellular geographic reference and boundary coordinate types of data bases contain the boundary information necessary for spatial display in mapping form. These information systems are referred to as cell and polygon systems. Image based information processing can also be included as a cellular type. These two primary systems are employed for their mapping capability.

#### Data Organization

The processing systems are markedly influenced by the data organization and can be characterized by the approach used to handle the geographic or spatial reference.

With the image raster approach the spatial or x-y position is implicitly recognized by the scan position within the image [10]. The processing approach is described as image manipulation [6]. A polygon data organization encloses each homogeneous spatial region with a polygon. The x-y coordinate pairs for the curved line segments which enclose the polygon are recorded.

The grid referenced or cellular approach divides the spatial domain into uniformly spaced and sized cells (a grid). The relative position of the cells in an array preserves spatial relationships without allocating storage to the location information [11]. Steiner [7] characterizes the grid based coding schemes as (1) complete - x and y coordinates and datum recorded for each geographic point, (2) sequential - only datum values recorded in a known sequence, and (3) compact sequential where string length or change point methods compress the sequence. String length coding involves a datum and a repeat factor while change point coding involves an initiator coordinate and a datum which applies until the next initiator coordinate. All three grid based schemes accomplish the same geographic referencing but processing economics differ.

The cellular and polygon approaches will be compared further. The image approach is the newest and least documented in comparative performance and costs. By virtue of its similarity to a full-matrix cellular approach, the image form will be dropped from discussion - realizing that continuing development in this area may warrant a separate review.

Steiner [7] cites grid system simplicity in data handling and output phases as advantageous. Bryant and Zobrist [10] compare the geographic encoding approaches in general terms as follows. The grid cell method is manually operated with poor spatial resolution and difficult update capability. The polygon method is expensive for large data sets and inherently prohibits certain operations. Continuing to summarize the drawbacks of both systems, they cite three primary disadvantages of each as follows:

- Grid cell - (1) spatial resolution tied to a cell
- (2) data are nominal or ordinal but not both
- (3) manual coding difficult/costly as is update process

- Polygon - (1) editing is computationally expensive
- (2) topological extraction of sub-areas is complex
- (3) computer system for overlay of polygons is expensive.

These summaries are indicative of a possible single tradeoff decision between cost and accuracy when choosing between these approaches.

To acquire a meaningful set of operating cost data for comparison of the systems is difficult because only one system is typically utilized and reported. Smith [12] overcame this problem by utilizing a benchmark project to be completed by contracting with two agencies - each representing one approach. Soil, slope and geology maps at 1:24,000 scale were provided for conversion to geographic data base and analyses. The grid system was operated on a 2.5 acre cell. To convert to data base, the costs were \$342 for the grid system and

\$2825 for the polygon. Similarly the analyses specified cost \$430 via the grid processor and \$4520 via the polygon system. The cell or grid system is much less expensive to operate.

The costs of data conversion versus analysis by the grid system in Smith's study were on the same order of magnitude. Tomlinson [5] reports that data preparation costs for cellular systems usually run four to five times the manipulation and display costs. Factors accounting for this apparent discrepancy are the quantity of analysis done once the data are prepared and a high variation in analysis costs according to the specific storage approach - complete, sequential or compact sequential - within cellular systems. The author has experienced cost ratios of 10-1 to as much as 20-1 in comparison of complete coding to compact sequential coding data base processing.

Operating cost factors appear to favor the grid system. There are criticisms of the grid system which warrant consideration. Smith [12] reports that grid type maps are criticized for quality of final product - especially where a line printer is used - and they are criticized for data precision in a grid framework. Both criticisms are answerable - the first by selection of another output media from among the ever increasing list of alternatives and both by utilization of small cells.

Phillips [9] and Sinton [13] provide a comparative summary. The cell system has grid coarseness as a limit on data resolution and the finer the grid the greater the storage requirements. It does work best and fastest on large or complex data sets because of data accessibility and adaptability to many cell-oriented devices, line

printers, film recorders, computer memory etc. The polygon system has excellent spatial definition and therefore a better map accuracy and product quality. The storage and retrieval phases are relatively more complex and difficult and the processing of data is more costly with the polygon system.

Sinton [13] notes that polygon users are typically cartographically oriented which requires the more precise ground locations while grid-cell users are resource analysts desiring manipulative ability and operating economy.

Smith [12] draws a conclusion in his study which well summarizes the comparative discussion here. The polygon system provides better quality products but is usually too expensive. The cell system is operable and affordable.

### Performance Evaluations

The consumer of information is concerned with the accuracy of the two primary analysis products, spatial inventories (tabulations) and spatial displays (maps).

System evaluation requires monitoring successful analyses for accuracy of inventory, accuracy of mapping and corresponding costs.

Switzer [14] outlined a mathematical approach to evaluation of the grid system. His evaluation was of map accuracy. Switzer's evaluation of map precision was based on discrepancy between an "estimated" map from the system and the "true" map (actual spatial relationships). In the set theoretic sense, the spatial intersection, defining the area of agreement, would measure map precision.



The coding method employed was the cell center datum. Whatever the cell size, the cell is considered homogeneous and of whatever data category or value occurred at the cell center. This differs from the cell dominant approach where spatial plurality within the cell determines assigned homogeneity. The cell size was recommended to be smaller than twice the smallest map region to be retained. This constraint arises obviously due to extreme error situations when many strata appear within one cell and central datum alone defines the result. When cell sizes are small so that map boundary lines on a cell basis look like linear segmentors, then there will be complete agreement of cell central datum coding and cell dominance coding techniques.

Switzer's mathematical formulation of the coding process as it relates to mapping error is diagrammed in Figure 1. A single cell is shown. Datum assignment would be to stratum  $j$  based on membership of central point  $B$ . In a mapping sense, an error arises due to omission of stratum  $i$  points from representation even though they are within the cell and a distance  $d$  away from  $B$ . The boundary line is drawn linear as would approximately be the case with small cells. In this situation cell dominance would have to agree with central datum in assignment of stratum  $j$ .

Geometric probability can be directly employed to evaluate  $P_{ij}(d)$ , the probability of stratum  $i$  points within a cell assigned stratum  $j$  and separated from the central datum by distance  $d$ . Clearly the area mis-match between the coded cell - entirely stratum  $j$  and the original "true" cell is exactly the proposed estimator of mapping precision. Thus, map precision depends on  $P_{ij}(d)$ .

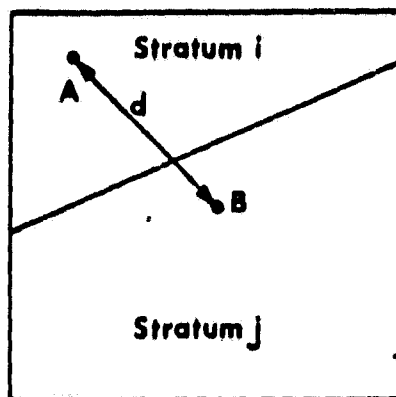


Figure 1. The central-datum cell coding concept as related to mapping error. Cell assigned to stratum j because of membership of central point B even though there is stratum i present in the cell at points such as A a distance d away from B.

Switzer evaluated cell geometries. For a cell of given size, he reasoned,  $P_{ij}(d)$  depends on the map and the cell geometry. Therefore, it is only necessary to hold the map constant as well as the cell size and vary the cell geometry. Resultant changes in map precision would provide the evaluation. Performing precisely this investigation he concluded that square cells are superior to rectangular and are only slightly inferior to a hexagonal network where all adjacent cells are equally spaced. Also, he noted that map precision or accuracy increases at the same rate as the maximum cell dimension decreases (increasing number of cells).

Switzer continued his study by holding the sampling geometry constant, i.e. square cell of given size. He reasoned that  $P_{ij}(d)$  is then a map property. The estimated map (after digitization) can be analyzed to predict the actual property  $P_{ij}(d=0)$  and that property in turn used to predict the map precision. This predictive sequence estimates map precision from the estimated map alone without a "true" map as reference. This is a commendable effort in mapping evaluation and the only assumption made is that map data are fairly dense.

The drawback to this evaluation is that the estimated map, i.e. the digitized representation, must exist before the map precision can be predicted. If the precision turns out unacceptable to the user, he is forced to redigitize to achieve a new estimated map. In operating characteristics noted earlier, cellular systems involve as much or several times more investment in the data entry than in analysis. Hence an iterative digitization and evaluation is not an economical approach to performance oriented digitization. What is needed is a map analysis method for performance prediction which can be employed before digitization begins.

Recent articles have suggested probabilistic models for the evaluation of match between a computer output map and reality. Recall this is exactly Switzer's approach to evaluation of map precision. Therefore, these models may also prove useful in evaluation of geographic information systems.

Van Genderen's view of land use map accuracy parallels Switzer's map precision. Cell by cell match of the computer map to actual ground conditions measure accuracy. He noted, however, that total comparison is often physically and economically impossible. A sampling technique is proposed. He suggests using stratified random sampling by category and a binomial probability model for the two situations, match and mismatch. His concern was how many samples,  $N$ , should be taken if under probability of error,  $p$ , one does not wish to risk a chance result of no errors in the sample. The objective of the sampling in the first place is to estimate  $p$ , but he is concerned with overconfidence due to misleading "perfect" scores.

Van Genderen's discussion is actually a specialized application of Hord's [15] earlier work along the same line. Hord proposed the binomial approach and expounded on use of the normal distribution to calculate confidence intervals on the estimate of  $p$ .

It should be noted that the Hord-Van Genderen probabilistic mismatch model, like Switzer's output map analysis, requires the product map be generated before evaluation can be made. This is not appropriate for judgements prior to digitization. It would be useful to analyze the map prior to digitization, to enable predictions of mapping and inventory accuracy, with various cell sizes.

#### Grid Cell Size Selection Bases

To access the analytical, manipulative and display capabilities of a geographic information system, the user must provide for conversion of his map or spatial data into the format of the system data base. Users of a grid cell information system face the obvious and extremely significant question of appropriate cell size. Crude guidelines do exist in the literature. Suggestions are as follows:

- (1) use the resolution of the source data [11,16]
- (2) consider processing capability and cost and use finest grid affordable [7]
- (3) use grid appropriate to the particular data application [16,17]
- (4) grid-cell size selected according to data detail such as urban at 10 meter, urbanizing at 250 meter, other at 25 kilometer [7]
- (5) grid-cell selected according to output device i.e. if line printer has rectangular characters use rectangular cells [7,2]

(6) grid cell size small enough to force smallest spatial unit on the map to be more than 50% of the cell [9,14]

(7) consider the volume of data generated for processing [11].

These guidelines for cell size selection suggest that knowledge of accuracy or performance capability of a grid system is required. There is a hint that a study of a cell size-accuracy relationship is needed.

With data conversion (digitization) under a grid oriented system being a significant part of the operating cost, the user might well hesitate to select a very fine cell spacing. Meyers, Durfee and Tucker [16] make the very valid point, however, that use of a smaller grid may generate more cells in the encoding process but may still save time by making dominant theme assignment trivial compared to large cells covering several themes. It is nevertheless recognized that a significant factor in cell size selection is the time and cost of the coding method used [11].

#### Research Objectives

All users of geographical information systems have a vested interest in learning more about appropriate application of the system. The effectiveness of a system can be judged in terms of application costs and product accuracies. With a grid cell system both of these measures depend on the cell size selected. The literature reviewed represents an effort to define a cell size selection and an effort to evaluate mapping and inventory products based on system/map characteristics.

The objective of the research herein reported is to overcome two serious shortcomings in the work to date - namely (1) that mapping and inventory accuracies have only been analyzed/modelled after the fact rather than prior to digitization of the data base and (2) that cell size selection criteria have not been based on the cost-accuracy tradeoff which appears to be present. The desired outcome of the research is a procedure for analyzing candidate map data for parameters which will assist the user in deciding on cell size.

A performance prediction procedure should hold for all cell-dominance coded, cellular geographic information systems regardless of the analyses they perform.

The route of investigation is to answer the following questions:

- (1) Do mapping/inventory product accuracies relate to cell size given a data set?
- (2) Do characteristics of the input map vary predictably with cell size?
- (3) Which characteristics are useful in accuracy prediction?
- (4) What procedure or technique might be employed to predict accuracy for a given map?

#### SYSTEM PERFORMANCE EXPERIMENTS

A research objective was to derive an accuracy prediction technique based on map characteristics and useable prior to digitization of the map. The route outlined for achievement of this objective included investigation of performance of the geoinformation system as

measured by inventory (tabulation) accuracy and by mapping (spatial) accuracy. The system employed was AREAS, Area Resource Analysis System. The performance is dependent on the grid-cell nature of the data and not on the specific programs of AREAS. Experiments on inventory and mapping accuracy as related to cell size variation are reviewed.

### The Data Base

A sample map segment of soil survey data was selected for intensive study. An area of moderate boundary density and a mixture of region sizes and shapes were selected. The intent was to avoid bias that might enter artificially created data or that might arise from data with regularly shaped and/or sized patterns. The data map segment is shown in Figure 2.

The utilization of a very fine grid network was realized as particularly important to the study of cell size influence. Initially cells smaller than the smallest separation of adjacent boundary lines were considered. In addition a constraint was imposed that the cell size be an even integer divisor of the commonly used 1.008, 4.032, and 16.128 ha (2.5, 10, 40 acre) cells. Under this constraint, the study of decreasing resolutions would pass through cell sizes which users recognize.

The final selection of a 0.007 ha (0.0174 acre) cell resulted in a grid of 384 elements square. The data set created was considered the "true" map as reference for all analyses. This map is shown in

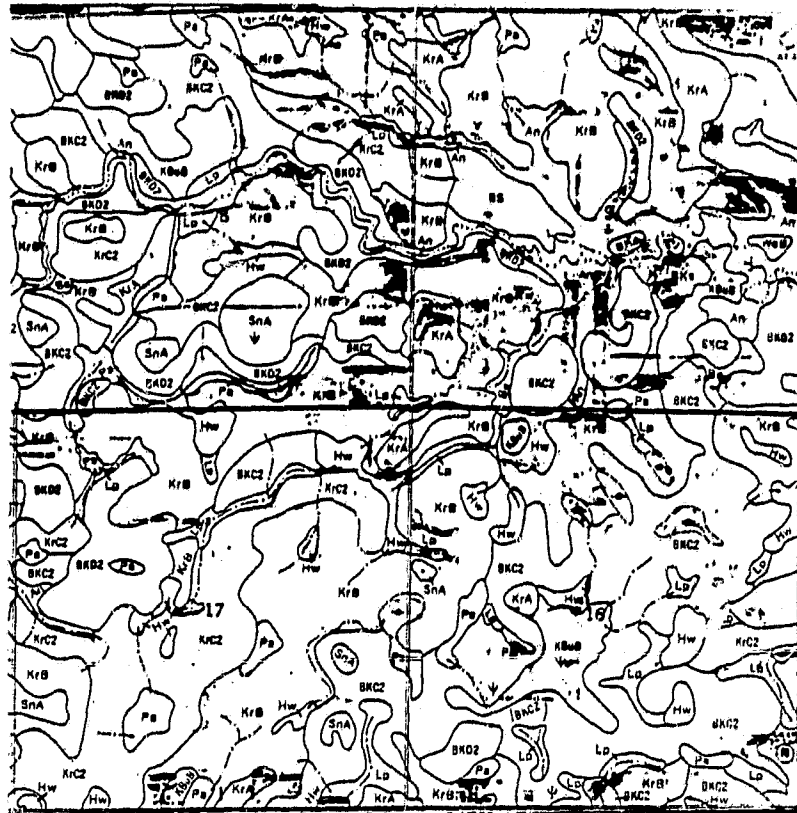


Figure 2. The original map of the intensive-study data set. The four sections of soil survey data are located in Minnehaha County, South Dakota.

Figure 3. The data set was generated by enlarging the map, preparing a computer drawn grid and manually encoding the cell contents.

For the analysis envisioned additional data sets of the same map at larger and larger cell sizes would be required. Rather than attempt many manual digitizations, the cheaper, faster and more dependable route of computer aggregation was chosen. An integer number is specified, for example four, and the change points of the data set are altered to align on column and row boundaries divisible by four. Four by four arrays of cells are examined to determine the



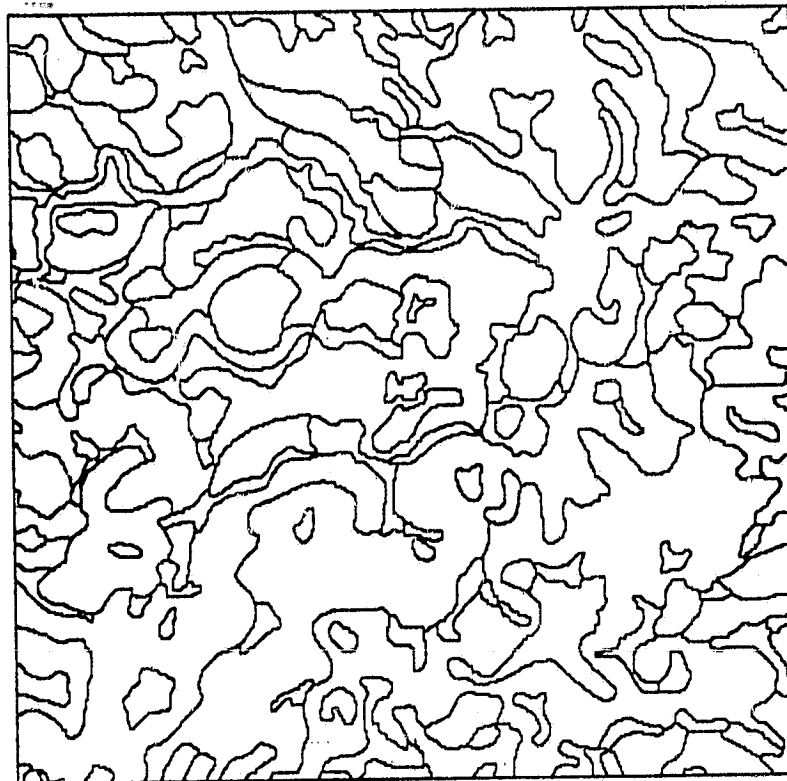


Figure 3. The computer map of the "true" data set at 0.007 ha (0.0174 acre) cellular grid.

dominant theme. Simple plurality applies and m-way ties are broken by an m sided coin toss - a random number process. The output data set size is maintained in agreement with the input data set.

For uneven division of the row-column dimensions of the input data set, the aggregation of pixels at the right and bottom edges would introduce artificial structure into analyses of boundary spacing. Therefore aggregation was only applied for even divisors of the 384 by 384 base data set. The aggregations were 2,3,4,6,8,12,16,24,32,48, and 64. Corresponding cell sizes were 0.028, 0.063, 0.112, 0.252,

0.448, 1.008, 1.792, 4.032, 7.168, 16.128, and 28.672 ha (0.069, 0.156, 0.278, 0.625, 1.111, 2.500, 4.444, 10.000, 17.778, 40.000, and 71.111 acres).

The eleven aggregations were also processed to reduce row-column dimensions. Rows were omitted and column designations divided by the aggregation factor. These forms of the aggregations were useful for more rapid and economical plotting as well as boundary structure analyses.

The original data set, the eleven aggregated data sets and the eleven aggregated-reduced data sets comprise the twenty three data sets of the data base for the intensive study.

#### Performance Experiments

Performance was analyzed from the product viewpoint as mapping accuracy and inventory accuracy. The influencing variable under study was cell size. The source data set with the finest resolution cell was the "true" map reference. An inventory tabulation of the spatial quantities of each mapping unit in the original data set became the "true" inventory reference for the study of inventory accuracy. The processing applied to each of the eleven data sets with increasing cell size is diagrammed in Figure 4. The inventory and mapping performance was observed by map units and summarized for the entire map. This allowed observation of size-shape influences as well as evaluation of the entire data set for each of the eleven different cell sizes.

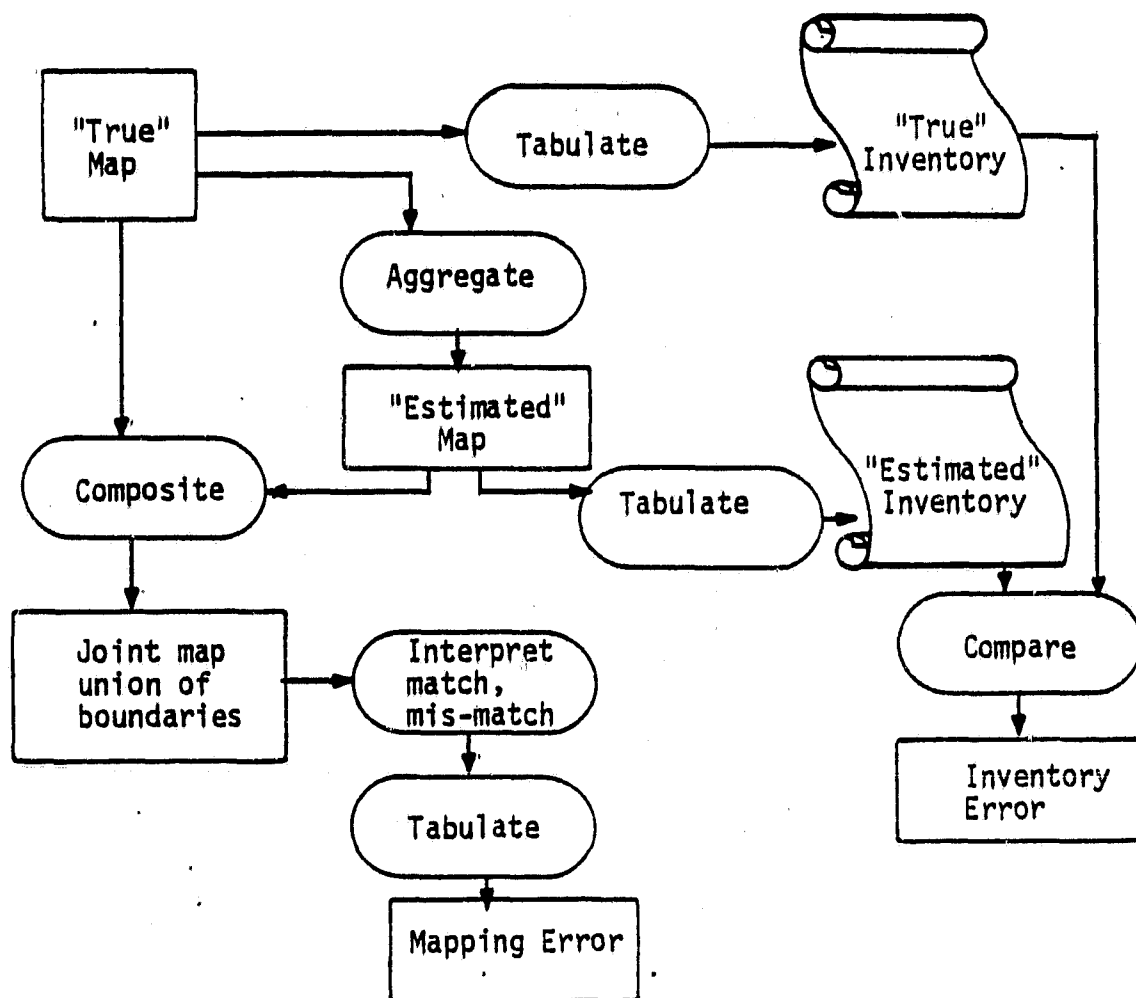


Figure 4. Performance evaluation processing diagram for an increased cell size.

### Experimental Results

The appearance of the maps corresponding to the twelve resolutions in the data base can be compared in Figure 5. Corresponding mapping errors, the mismatch of areas when comparing maps of larger cell sizes to the smallest cell size available, are displayed spatially in Figure 6. Note that the basic data set (upper left in Figure 5) and the aggregation by 2 to the immediate right are of such a fine cell size for the output scale being used that the cellular nature is not

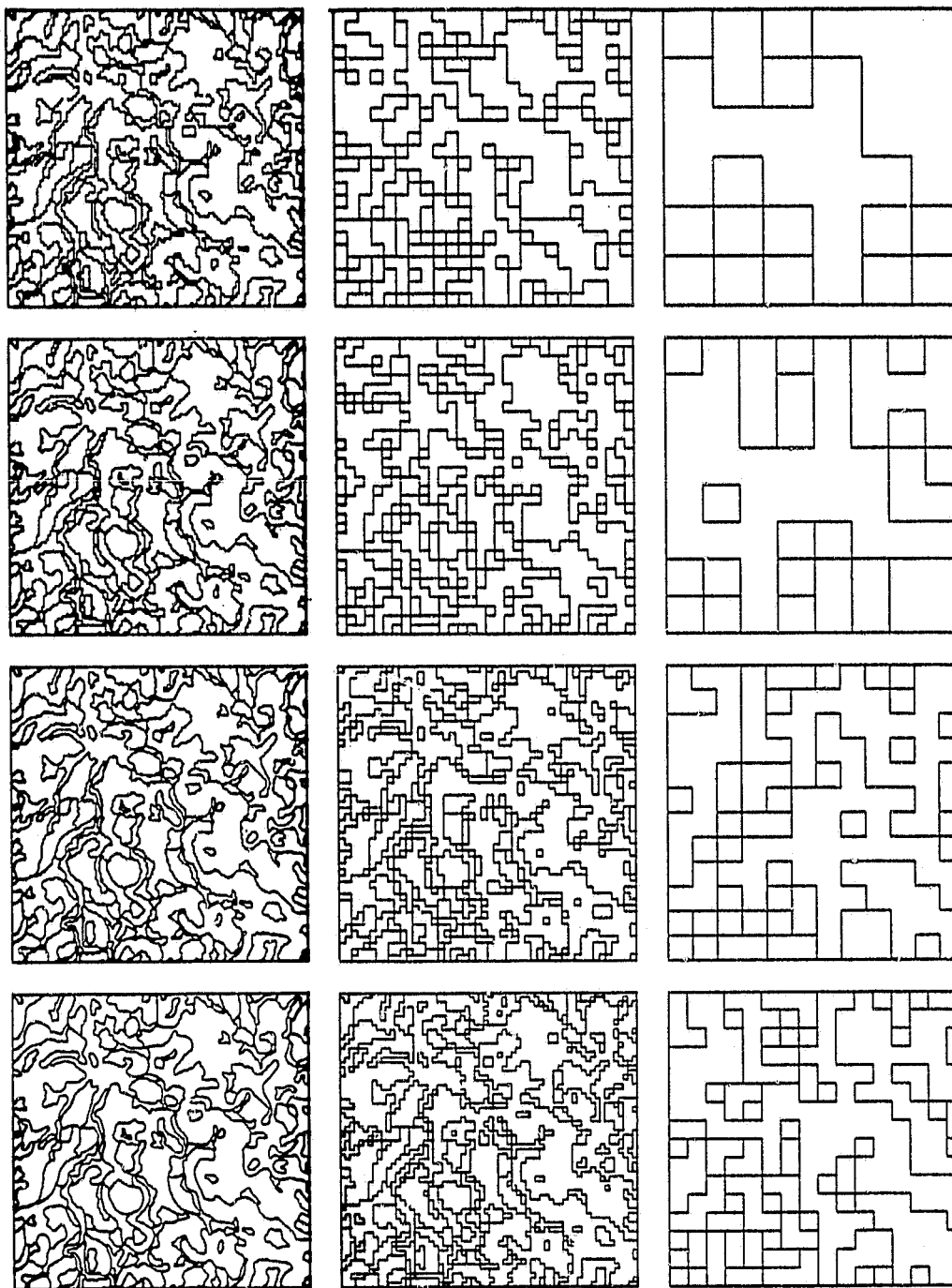


Figure 5. The twelve maps compared for mapping accuracy versus cell size. Cell sizes from top left to lower right are 0.007, 0.028, 0.063, 0.112, 0.252, 0.448, 1.008, 1.792, 4.032, 7.168, 16.128, and 28.672 hectares.

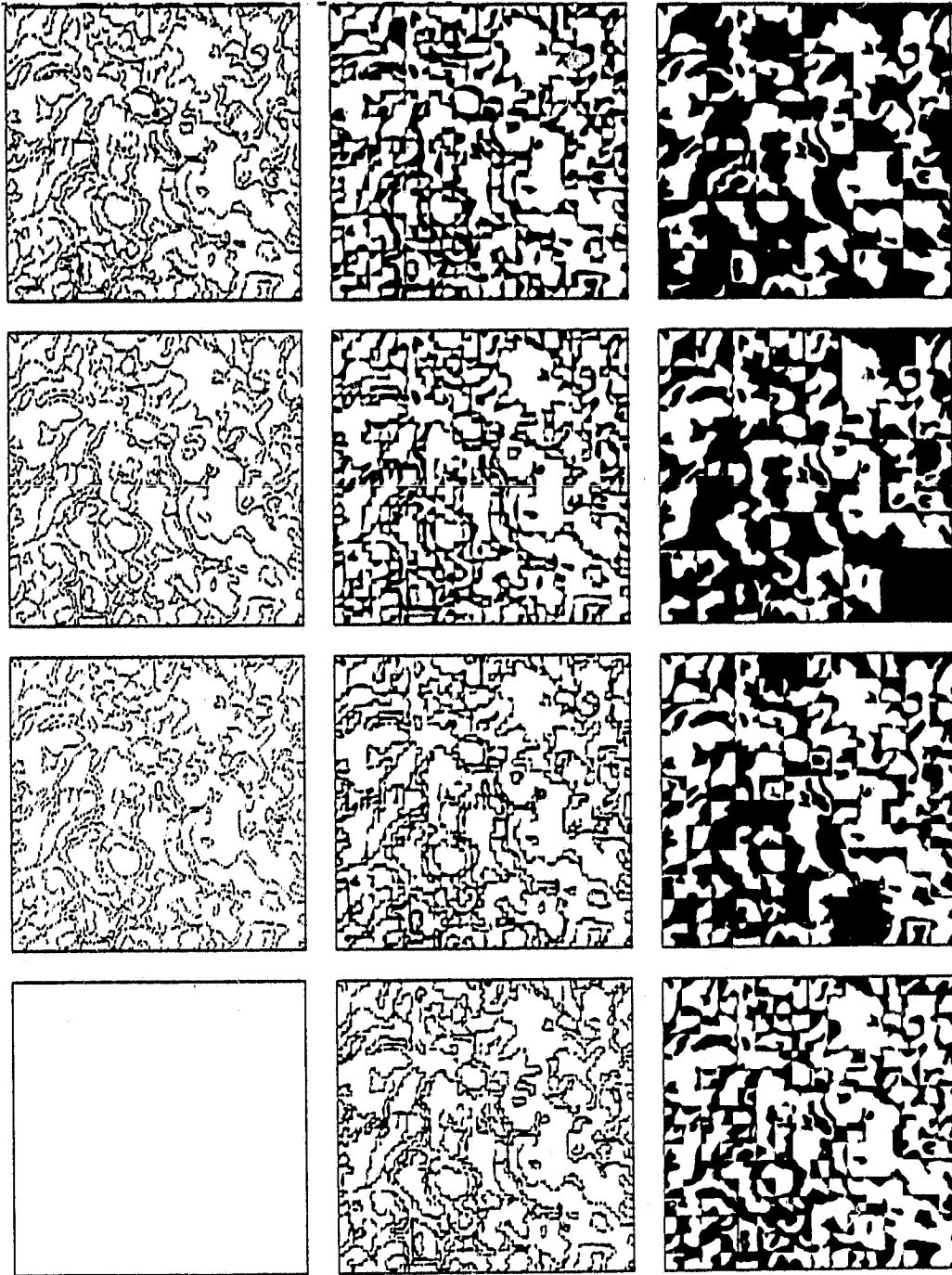


Figure 6. Representations of mapping error. Maps correspond to Figure 5.

ORIGINAL PAGE IS  
OF POOR QUALITY

even apparent. Cellular representations can be as cosmetically pleasing as the polygon approach if the fine cell size can be afforded.

Throughout remaining discussion of the data sets of different cell sizes a resolution number will be reported which is the aggregation factor for combining input cells from the basic map at 0.007 ha (0.0174 acres) cell size. Resolution numbers are 2,3,4,6,8,12,16,24, 32,48 and 64.

The twelve data sets represented in Figure 5 and the eleven non-zero data sets in Figure 6 all stem from map one (upper left) of Figure 5. This map has eighteen levels of soil association and all other data sets also had eighteen or fewer levels depending on the loss of detail.

When each of the coarser resolution data sets were composited with the first the resulting eleven, mapping-error data sets also contained multiple categorical match mis-match data. An interpretation process produced binary match-mis-match representations as shown in Figure 6.

The data sets corresponding to the figures 5 and 6 were tabulated. The tabulation of data corresponding to Figure 5 allowed numerical evaluation of the inventory accuracy versus the resolution number (cell size). The tabulation of data corresponding to Figure 6 allowed numerical evaluation of the mapping accuracy versus the resolution number (see Figure 4). The results obtained are plotted in Figure 7 versus the resolution number. The percentage mapping error is the number of incorrectly mapped cells divided by the total 147,456 cells

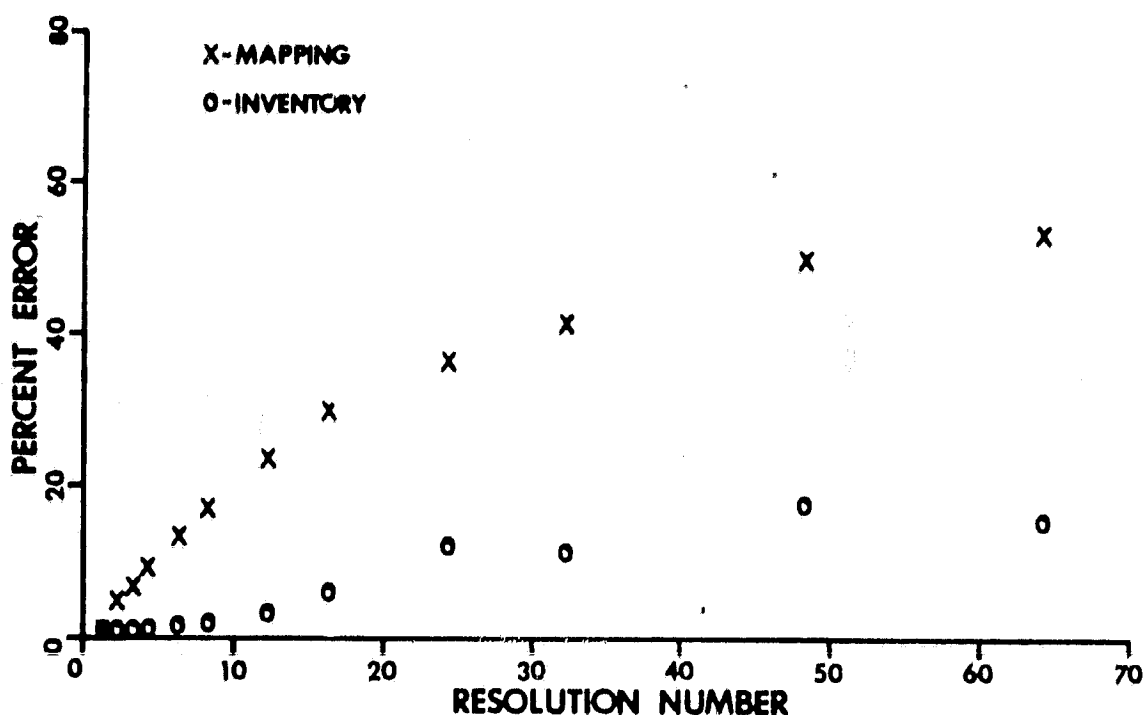


Figure 7. Mapping and inventory errors versus resolution number.

in the 384 cell square. Tabulation error will be lower than mapping error since omission-commission errors per category can offset each other. Since all cells are tabulated as belonging to some category, the average error across categories, with some categories over tabulated (+ error) and some under tabulated (- error), will always be zero. Root-mean-square, RMS, error was considered to overcome cancellation but the value obtained would be an average value for a map category - not the entire map. A root sum square, RSS, was employed where the squares of category tabulation errors are added and the square root taken. The result is an accuracy figure for the entire map. These are the percentage inventory errors graphed in Figure 7.

### Performance Relationships

The evaluation processing diagrammed in Figure 4 generated data sets which allowed performance comparisons by mapping unit or category within the map as well as by total map. The eighteen categorical units in the map were represented by 172 diversely shaped and sized spatial regions.

Shape and size are logically the characteristics of the mapping units which interact with any selected cell size to give rise to the mapping and inventory errors under study. These characteristics relate directly to the adequacy of representation by uniform cells. Larger mapping units will accommodate larger cells. Also larger mapping units will tend to have fewer perimeter cells in relation to total cells. Shapes may be broadly classified into a complexity spectrum from simple to complex. The criterion is border versus area or perimeter cells versus total cells. Simple shapes will require less perimeter to enclose an area and, therefore, tend to be less subject to the errors which arise in cellular representation of the borders. Figures 8 and 9 show the behavior of mapping and inventory errors versus resolution number for several selected mapping units of differing size and shape. These mapping units are separated and spatially displayed in Figure 10. Table 1 summarizes characteristics of these mapping units as a reference for studying the error patterns of Figures 8 and 9.

The anticipated trends in error behavior appear to be present. The lower error rates with changing cell size are associated with the larger mapping units. The small, single-region mapping units are



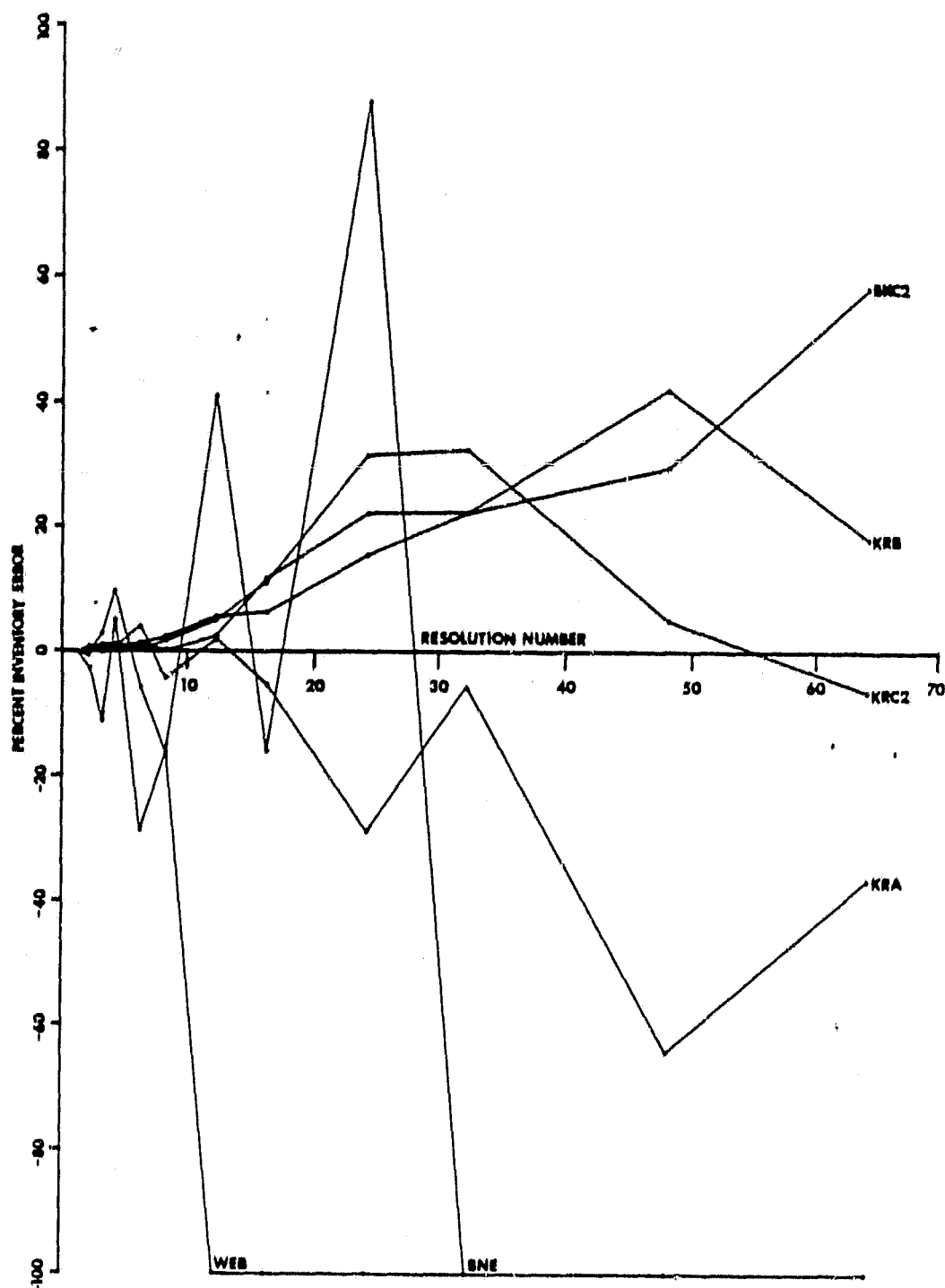


Figure 8. Inventory errors versus resolution number for selected mapping units.

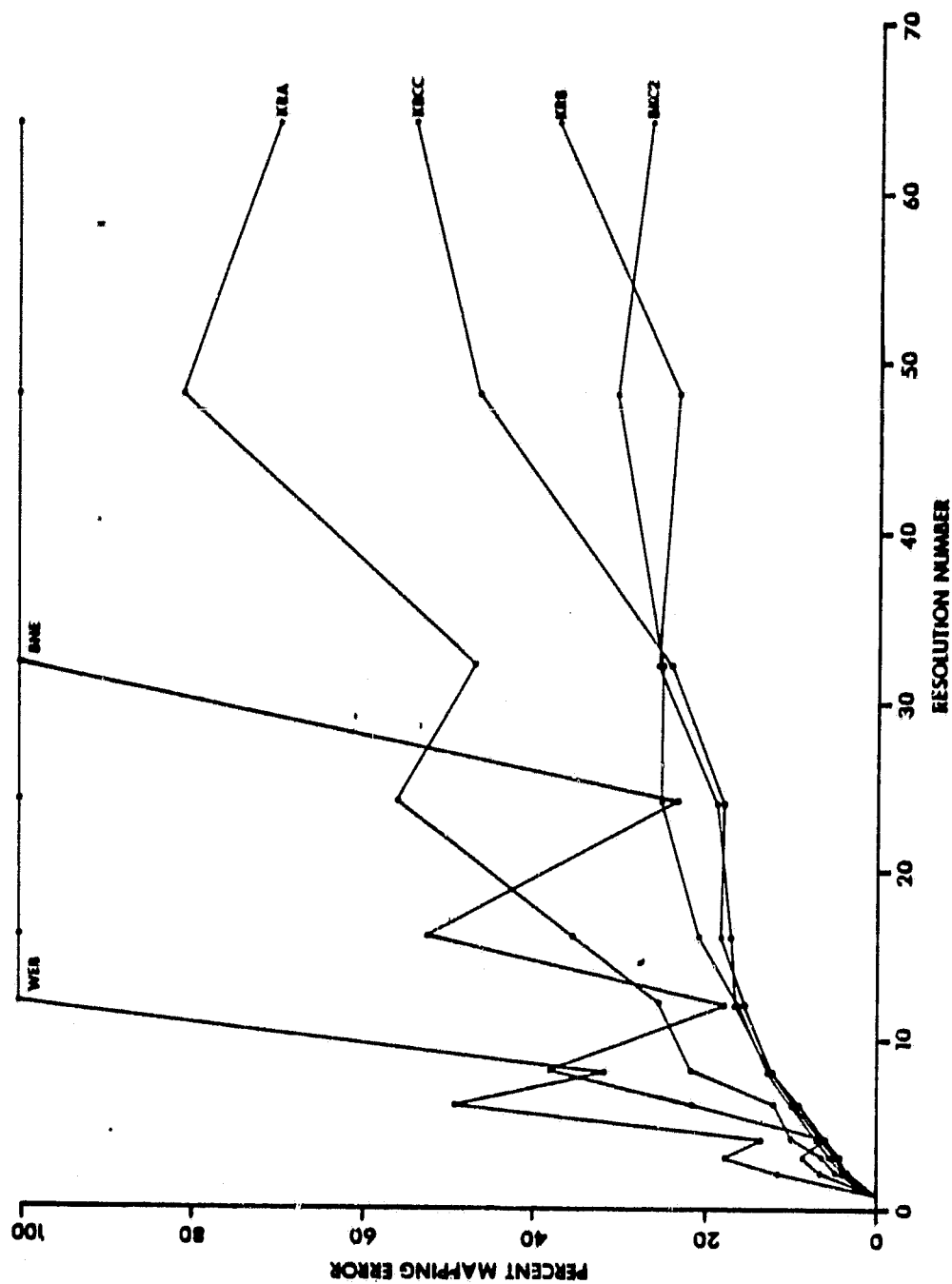


Figure 9. Mapping errors versus resolution number for selected mapping units.

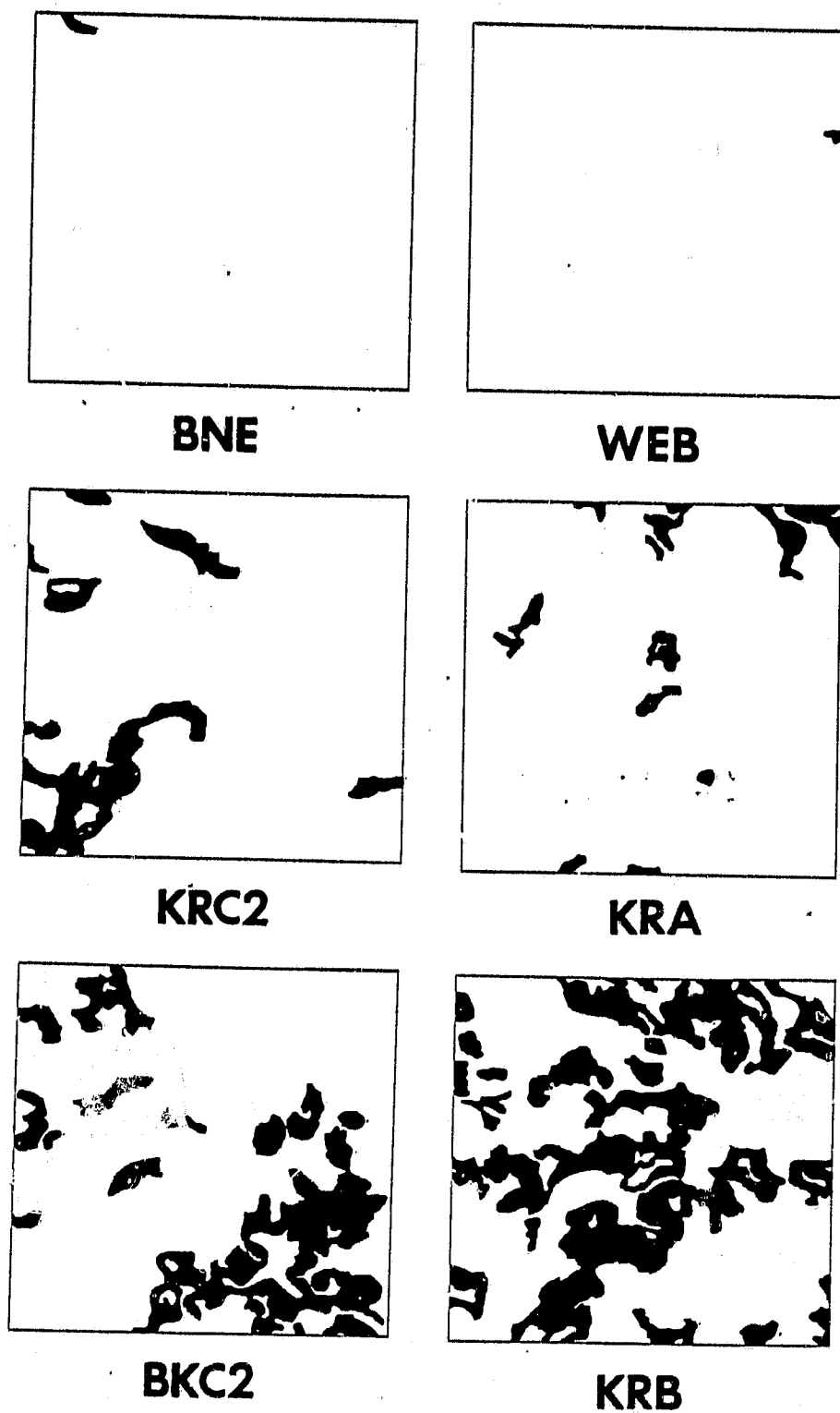


Figure 10. Spatial representations of the mapping units referred to in Figure 8.

Table 1. Mapping unit characteristics for mapping units referenced in Figures 8 and 9.

Mapping Unit	No. cells	hectares (ha)	acres (a)	No. regions
WEB	153	1.071	2.662	1
BNE	307	2.149	5.342	1
KRA	6530	45.71	113.62	11
KRC2	13199	92.39	229.66	7
BKC2	28555	199.89	496.86	10
KRB	48778	341.45	848.74	18

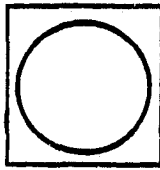
subject to an oscillatory relationship between error and changing cell size depending upon the spatial relationship of any particular cell size to the mapping unit. Once the cell size is large enough to assure that no orientation of that cell, with respect to the mapping unit, will retain any representation of the mapping unit, then the error goes to 100% and remains constant.

Spatial area alone does not totally determine the error versus changing cell size. The largest three mapping units are not strictly in the same performance order as size when comparing mapping and inventory behavior. The additional influence of perimeter is evident as partially determined by the number of regions of the mapping unit. Shape of the mapping units also determines the extensiveness of perimeters in relation to area.

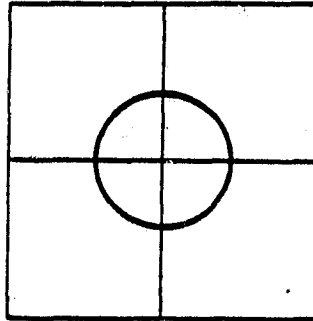
#### An Orientation Study

Beyond considerations of the mapping unit shape, size, area and perimeter, there is also the potential impact of grid position.

Figure 11 demonstrates the extremes in mapping error which might be



Mapping error = 0%  
(a)



Mapping error = 100%  
(b)

Figure 11. The effect of grid position on mapping error for a single circular mapping unit. In part a the circle dominates the cell whereas in part b it fails to dominate any of the grid cells. Mapping error is with respect to the mapping unit.

observed with two alternate positions of a particular grid cell size. It is evident that isolated mapping units may exhibit widely varying mapping errors dependent on grid position.

In the study of cell size effects, various cell sizes were generated by aggregation of cells in the reference map. The grouping of  $k$  by  $k$  cells into a larger cell was accomplished for all integer  $k$  which evenly divided the map dimensions. Every aggregation regardless of  $k$  began with row one, column one of the reference data set. Hence each aggregation used in the study was one of many possible orientations for that aggregation. In fact for an aggregation factor  $k$  there are exactly  $k^2$  ways of spatially positioning the enlarged cellular network over the original map.

It was felt that some of the variations in error present in Figures 8 and 9, particularly for the smaller mapping units WEB, BNE and KRA, arose from the particular positional orientation of the aggregated cells. To evaluate the potential magnitude of this variation, to seek any underlying average trend, and to gain insight into data structure relationship to error, an experiment was undertaken. Several simple closed figures were drawn on cellular reference background. A circle, square, ellipse and rectangle were used. The surrounding background was not considered of mapping interest and not included in mapping error calculations. For each integer aggregation  $k$ , the  $k^2$  ways of aggregating were all observed and mapping error for each calculated. The  $k^2$  observations of each aggregation  $k$  allowed calculation of a positional average mapping error. Errors arising from aggregation with respect to one common coordinate point were recorded separately to enable comparison of the average to the behavior observed with the procedure as used on the intensive study area.

Figure 12 contains the average and common-reference error graphs versus resolution number for the four simple regions analyzed. The suspicions regarding the presence of chance positional effects in Figures 8 and 9 appears justified. The error rate for average positional aggregation is a much more well behaved and non-decreasing function. In fact the proximity to a linear area-based model is quite pleasing.

The linear, area-based model simply stated is zero error at the reference or base cell size and linear increase in error rate to 100%

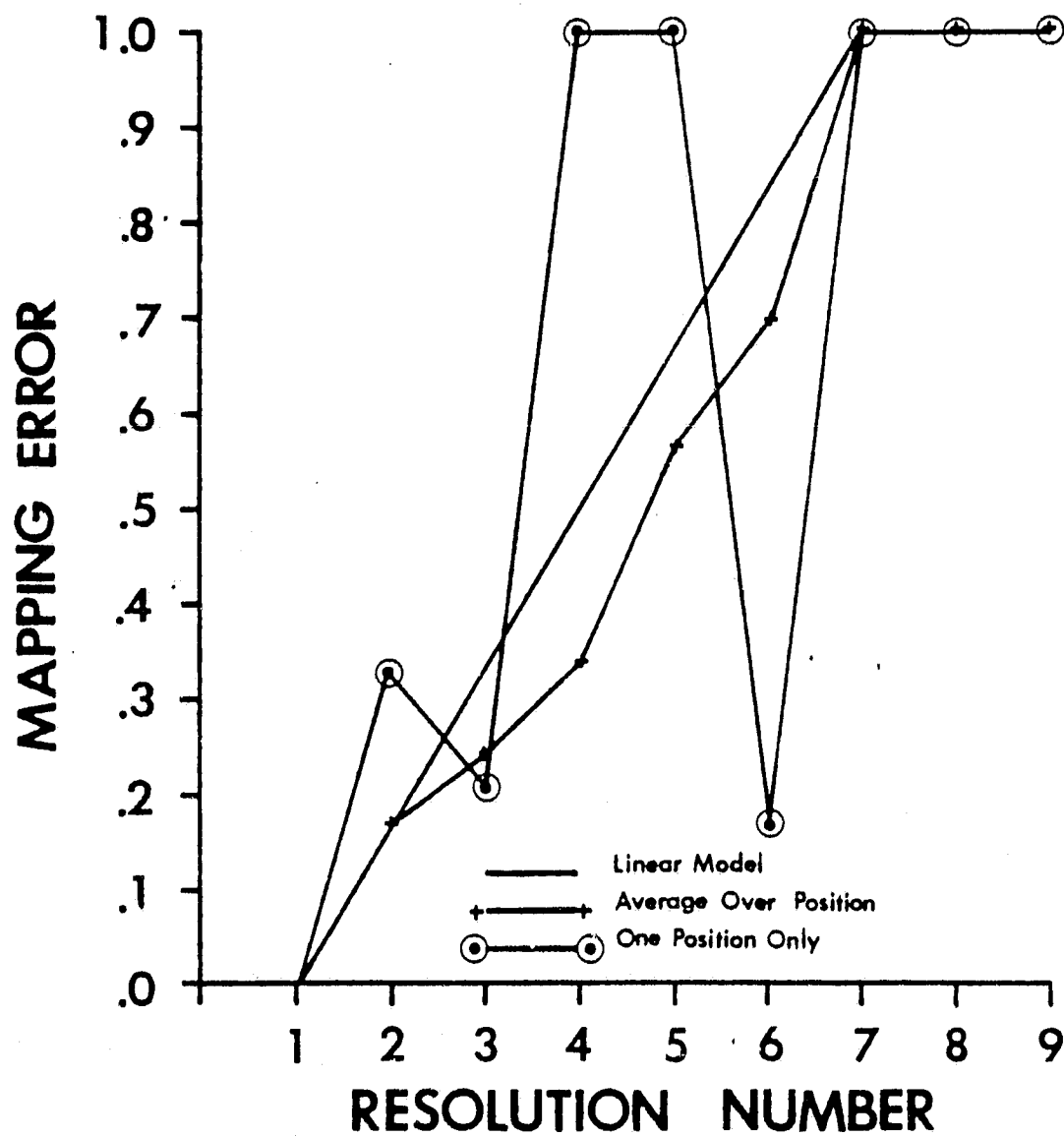


Figure 12 .a. Mapping error for various resolution numbers in the case of a simple closed circle.

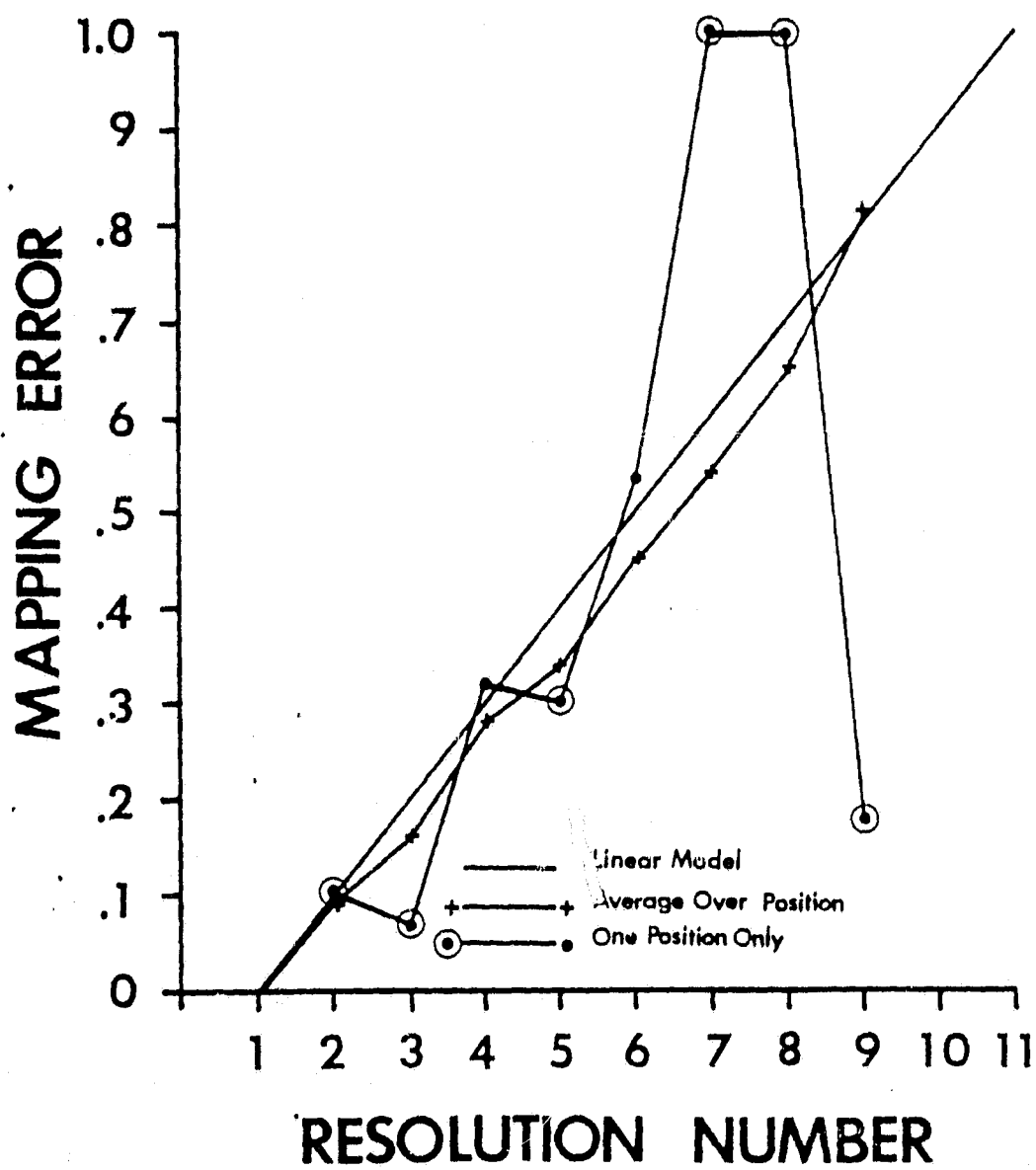


Figure 12.b. Mapping error for various resolution numbers in the case of a simple closed ellipse.



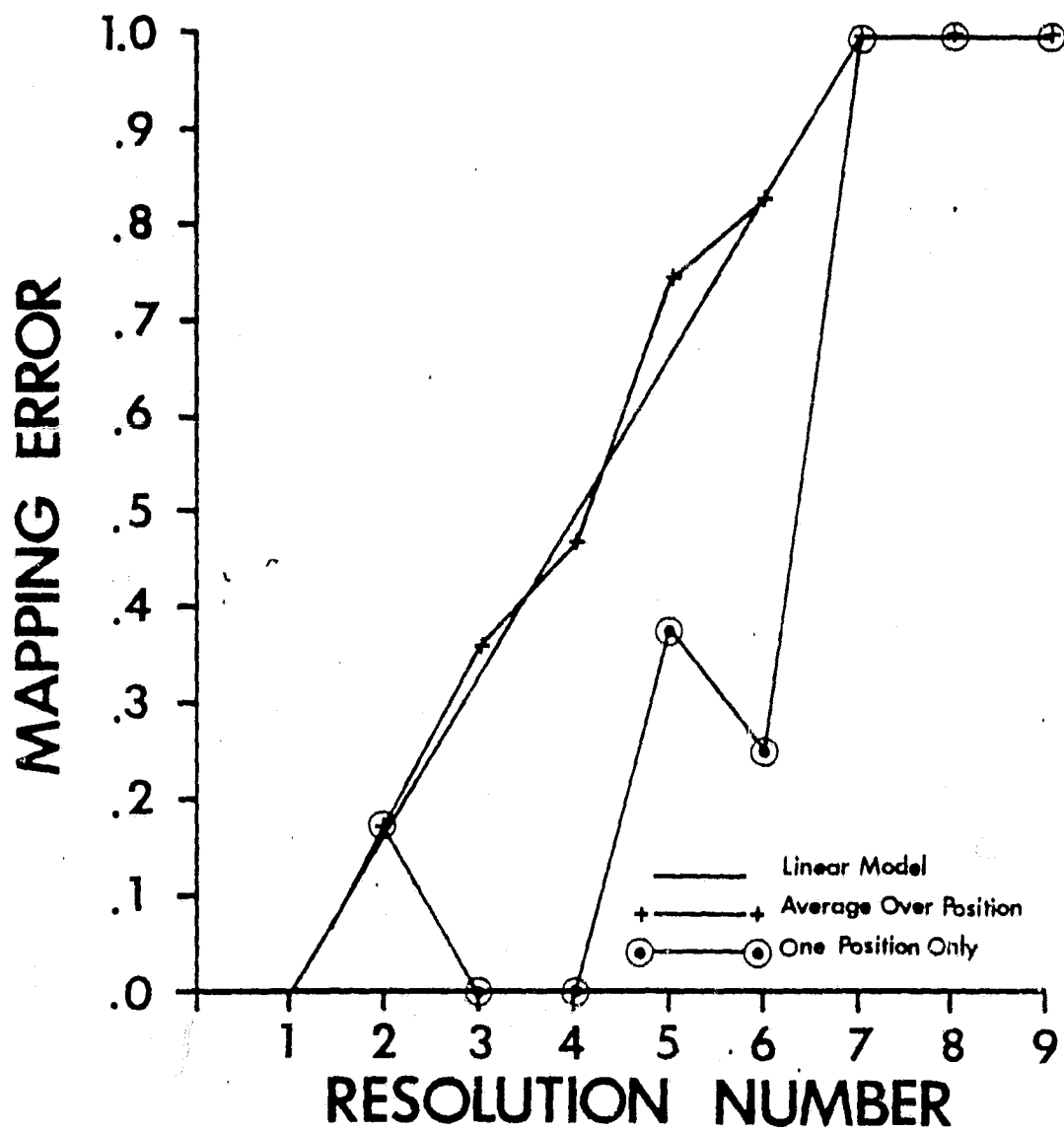


Figure 12.c. Mapping error for various resolution numbers in the case of a simple closed rectangle.

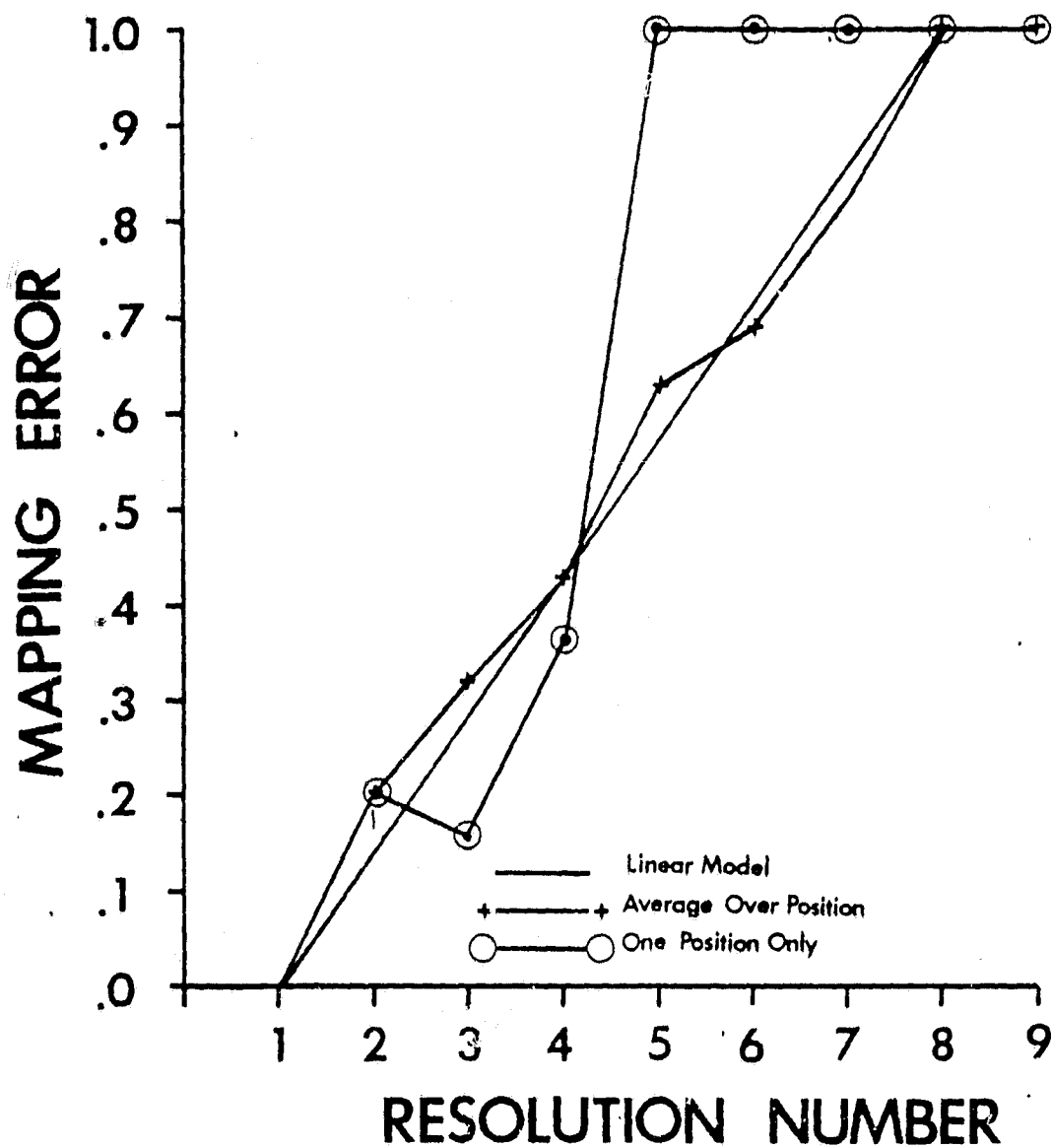


Figure 12.d. Mapping error for various resolution numbers in the case of a simple closed square.

at whatever cell size corresponds to a cell area greater than twice the area of the map region. The latter situation for a simple closed region occurs when the cell, regardless of position, is no longer dominated by the region. This is a simple, intuitively appealing and reasonably accurate representation of error relationship to changing cell size for simple closed regions isolated on a spatial background.

Reconsider the linear, area-based model for a rectangle 2 units by 55 units. The area of 110 square units would predict 100% error for resolution numbers greater than 14 ( $15^2 \geq 110 \times 2$ ). In reality, however, the dimension of 2 units would disappear for resolution numbers greater than 5. Clearly the linear area-based model applies only to restricted cases.

It is asserted as self-evident that the averaging over the many possible aggregation positions of cells with respect to a single closed region is equivalent to the averaging over many randomly dispersed identical regions with respect to a single aggregation position. This extends the experiment results to the case of the intensive study map and allows the experimental results to be used to explain the error behavior in Figures 7 through 9. The conclusion is also drawn that area is ambiguous in prediction of mapping error.

The conclusion is also drawn that mapping error relates well to changing cell size if effects of positioning the cell can be averaged.

This simply implies that analysis of mapping error versus cell size is justified over a map segment as observed in Figure 7 but not for isolated mapping units as in Figures 8 and 9. Even though several factors appear to influence the errors on a mapping-unit basis, the relationships of Figure 7, effectively representing averages over the eighteen map units, are adequately smooth trends which could be modelled by least squares curve fitting.

#### DATA CHARACTERIZATION ANALYSIS

The behavior of mapping and inventory accuracies as cell size changed was observed. The map characteristic/s relating to changing cell size were also analyzed.

##### The Distribution of Spans

The change-point version of the compact, sequential coding scheme is utilized by AREAS. The distance equivalent to  $n$  uniform cells is a discrete representation of the continuous interboundary separation in the input map.

From an intuitive standpoint, the distances which separate map boundaries are precisely the "character" of an input map to which cell size must relate in influencing accuracy of mapping and accuracy of tabulation. If all inter-boundary distances were much larger than the selected cell size, many cells would occur between boundaries and remain accurately encoded even though the boundary cells are assigned by dominance considerations. In other words, the area in error within

any given boundary cell would be a small portion of the area represented by the remaining interior cells of the same map unit. As cell size grows larger relative to a given map fewer cells occur within a region and the border cells of the region include larger area errors. This intuitive concept is diagrammed in Figure 13. The actual occurrence of this effect is quite evident in Figures 5 and 6.

Inter-boundary distances may take any value in the measurement continuum. In the cellularized representation only some integer multiple of the cell dimension may occur. Thus the distances between boundaries in terms of cells are discrete integers and these distances will be called spans. Spans obviously relate to the coding scheme in the row or horizontal direction as shown in Figure 13. Spans also arise implicitly in the column or vertical direction and relate to error in the same manner as horizontal spans.

There is only one unique case of input map type which would give rise to a single value span. That is a map of boundaries which are a uniformly spaced, square grid. Boundary maps of natural phenomena (soil associations, lakes, physiographic regions etc.) and management/political units (counties, watersheds, national borders etc.) are most often irregularly shaped and sized regions. Even in the case of land use where field shapes are typically rectangular, the sizes and orientations vary and center pivot irrigation gives rise to some circular patterns. Typical maps will have a multitude of span values occurring with differing relative frequencies.

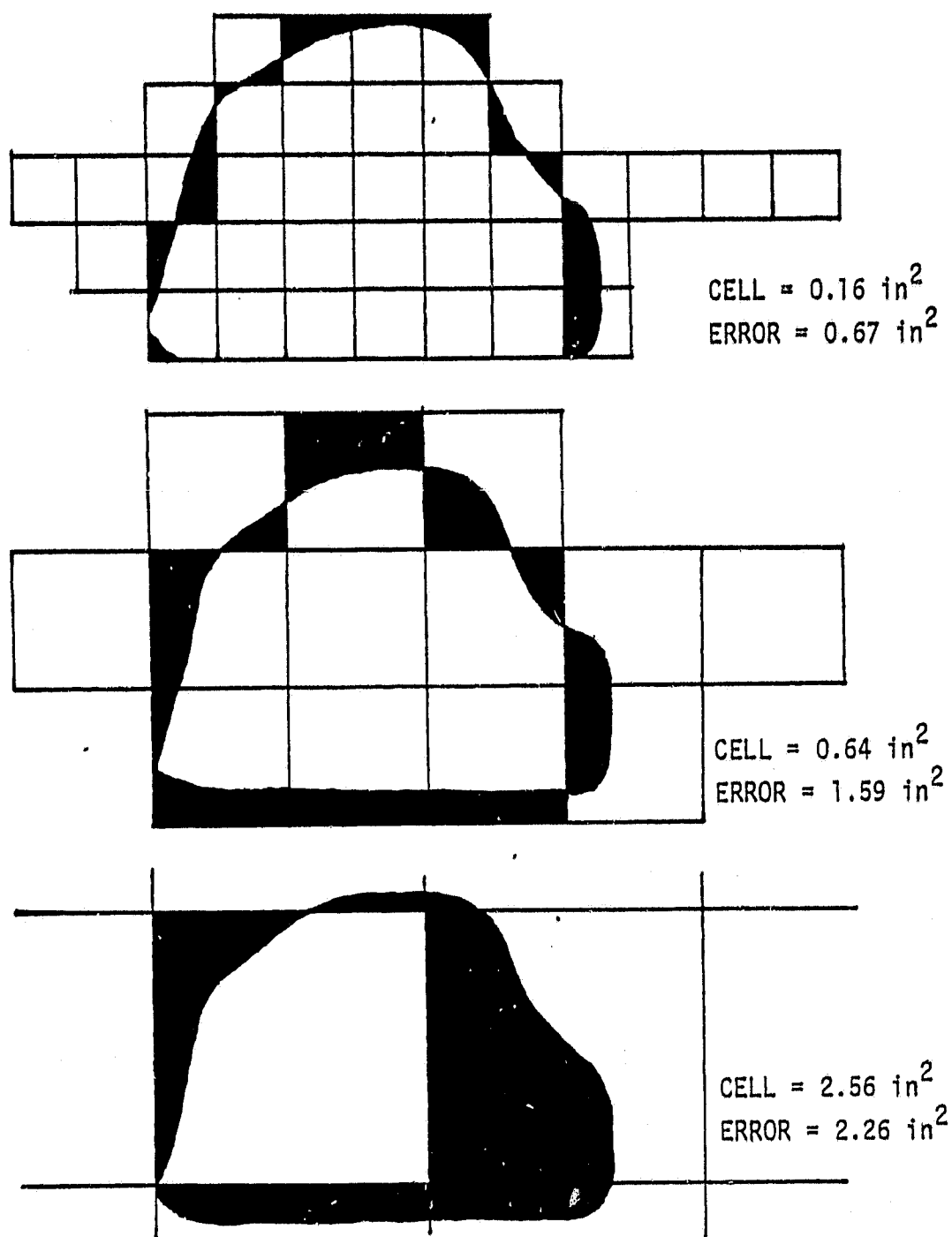
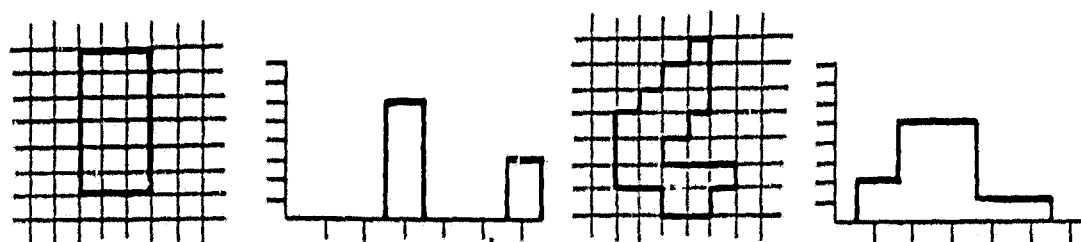


Figure 13. Growth of mapping error as cell size increases for fixed interboundary map distances. Area in error by cell-dominant coding process is shaded.

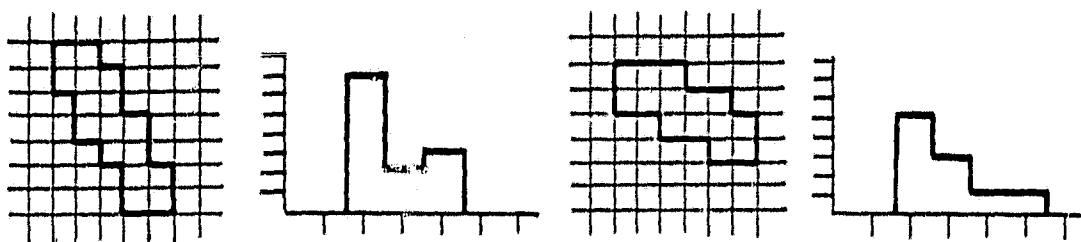
The influence of size and shape factors recognized earlier in relation to errors would be represented in the span distribution. The effect of orientations of an object would also be present in a distribution of spans. Figure 14 shows the effects of shape, orientation and size of a region on the span distribution. Note that the span distribution is shown as discrete in terms of the number of unit cells in the span.

An AREAS algorithm calculates span distributions and distribution means, and draws graphs of the relative frequency distributions for spans in the horizontal (row), vertical (column) and total (map) context. The Kolmogorov-Smirnov test is applied to the horizontal and vertical distributions to test the hypothesis that these sample distributions came from the same population for the class of data being analyzed. Rejection of the hypothesis in the statistical sense is conceptually an admission of directional sensitivity in the data set. The span distributions along and across rows of the data set are aligned with the cellularization in exactly the manner that the cellularization must impact inter-boundary distances. Thus this analysis provides characterizing distributions for any given cellularized data set.

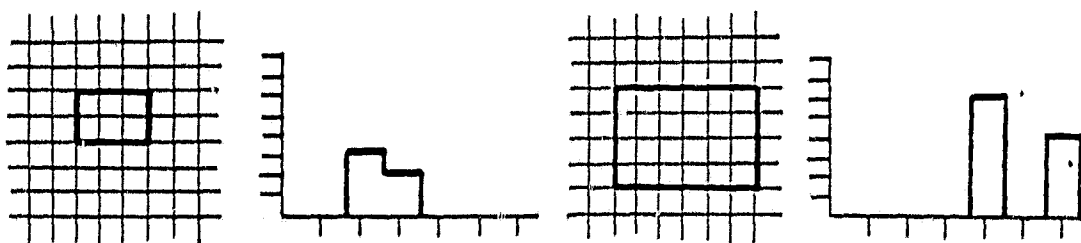
The finest resolution (smallest cell) data set, referred to as resolution one, was the "true" reference map. The three span distributions are shown in Figure 15. The Kolmogorov-Smirnov statistic was 0.052 with a rejection  $\alpha$ -level of 0.0000 to four decimal places. Rejection of the hypothesis of common parent population can be made with essentially zero chance of committing a Type I error, i.e.



(a) regular vs irregular shapes



(b) orientation factor



(c) size factor

Figure 14. The span distribution reflects region shape, orientation and size.



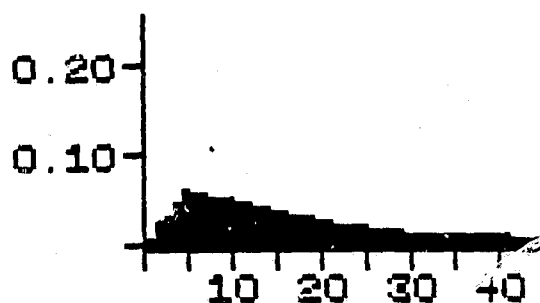
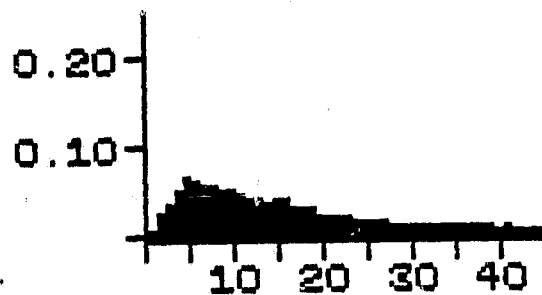
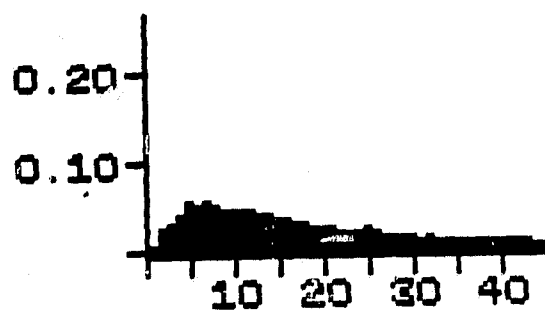


Figure 15. Span distributions for the reference data set of resolution number one. (a) horizontal scan, (b) vertical scan and (c) total map.

rejecting when in fact true. Although the distributions appear very similar, the statistical conclusion is that they are not.

This seemingly paradoxical situation is resolvable by considering the degrees of freedom for the test. The horizontal distribution had 7981 span samples for total map coverage and the vertical had 8748 span samples. With so many samples available the distributions would have to be nearly identical to statistically fail to reject the hypothesis of a common parent population.

It would appear likely that all maps except a rare class of very-highly structured and uniquely oriented maps would be spanned differently in the horizontal and vertical directions under small-cell analysis. Each of the study data sets, referred to by the resolution number, were analyzed. The total span distributions are shown in Figure 16.

The hypothesis of common parent population for the horizontal and vertical components of span was tested by the Kolmogorov-Smirnov statistic for the twelve resolutions with the results as tabulated in Table 2. Note particularly that statistical dissimilarity disappears at the resolution of six which is the mode of the reference data set. Apparently cellularization itself is beginning to cause the span structure in the horizontal and vertical directions to become statistically similar. This is suggestive of a test for too coarse a sampling cell size. The sampling of a map and generation of statistically similar horizontal and vertical span distributions is suggestive of a too-large sampling cell. The natural data structure should not be dominated by the cellularization. This is not to say

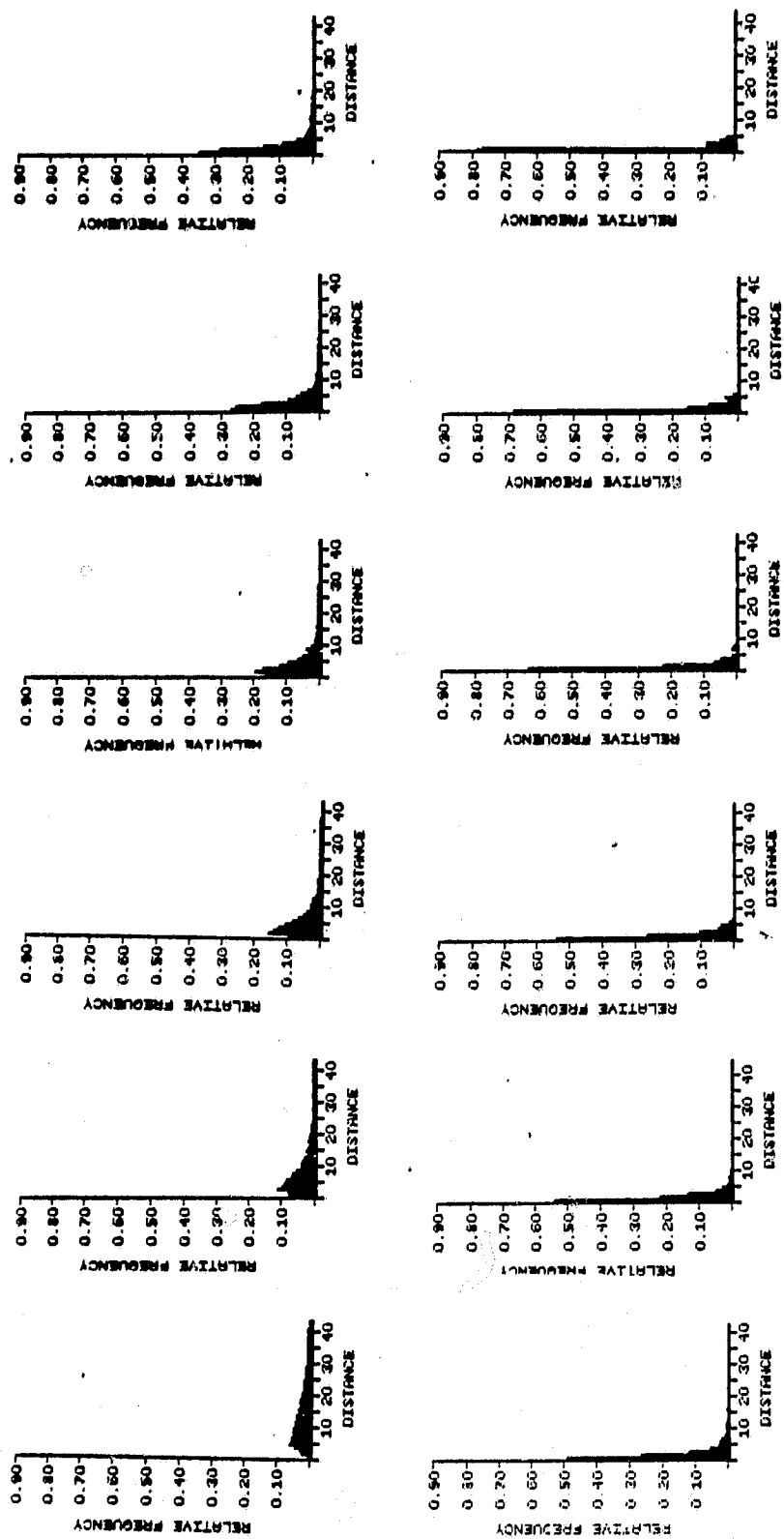


Figure 16. Span distributions for combined horizontal and vertical scans of the twelve different cell-sized data sets. Resolution numbers are (a) 1, (b) 2, (c) 3, (d) 4, (e) 6, (f) 8, (g) 12, (h) 16, (i) 24, (j) 32, (k) 48, and (l) 64.

Table 2. The Kolmogorov-Smirnov test of common parent span distribution for the horizontal and vertical span distributions of the twelve data sets.

Resolution Number	Kolmogorov-Smirnov statistic	$\alpha$ level	Conclusion on span distributions *
1	0.052	0.0000	different
2	0.046	0.0003	different
3	0.057	0.0003	different
4	0.065	0.0004	different
6	0.053	0.0631	same
8	0.042	0.9611	same
12	0.044	0.9408	same
16	0.042	0.9611	same
24	0.066	0.9215	same
32	0.071	0.9823	same
48	0.219	0.3221	same
64	0.076	1.0000	same

\* Hypothesis test at  $\alpha = 0.05$

that errors arising from cell sizes larger than the natural mode of the data are unacceptable. This remains an application judgement considering results such as in Figure 7 where tabulation errors never exceeded 20% even to the extreme aggregation to resolution number 64.

The key to success in the characterization analyses is whether the distribution selected (particularly a parameter of the distribution) relates to changing cell size. The trend in the distributions of Figure 16 is evident and to a degree is predictable. If the cells of all resolution one spans could be grouped by twos without relocating any boundaries in the aggregation process, then the span distribution for resolution number two would be that for resolution one with the span axis labels halved. Achievement of this in practice would require that all resolution one spans be divisible by two and that horizontal and vertical spans not be interrelated. This will not likely be the case. Hence the relationship of span distributions under changing cell size is not absolutely predictable.

To monitor and analyze the effect of changing cell size on the span distribution a parameter of the distribution had to be selected. The mode of the "true" reference data set was interestingly related to detectable difference in horizontal and vertical structure under changing cell size. The modes of the aggregated data sets, however, converge to one for the aggregation which surpasses the mode of the reference data set. Larger resolution number data sets have modes of one and discriminatory power via a mode parameter is lost.

The mean, however, varies with resolution number throughout and was selected for analysis.

#### Spans versus cell size

Each span distribution has a calculable mean at the particular resolution involved. The sample mean span was calculated from the estimated total span density function (as represented by the relative frequency distribution) by the conventional discrete formula

$$\bar{x} = (1/n)\sum x$$

recognizing that without knowledge of the specific form of the density function  $f(x)$  there could be no claim regarding unbiased maximum likelihood estimation. The calculated mean spans for each resolution number are plotted in Figure 17.

The relationship is logical. At resolution number one the span mean is the actual mean of the spatial structure. Map resolution one is a lower limit on the x-axis. As resolution number increases

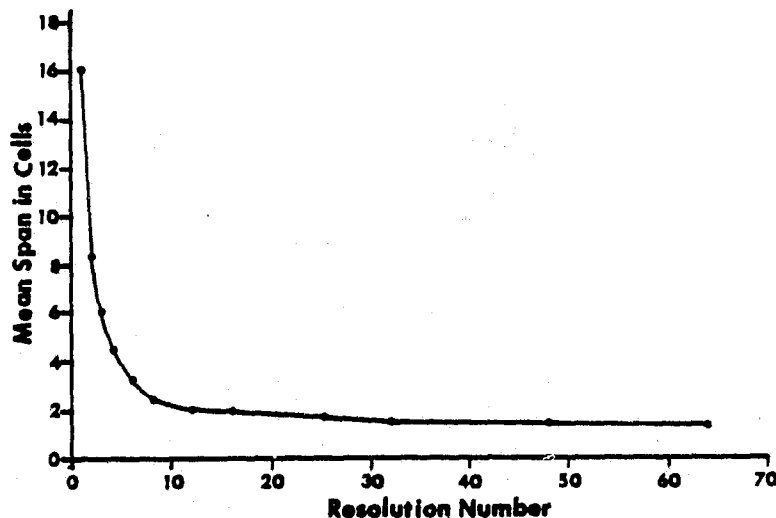


Figure 17. Span means versus resolution numbers

even to the full dimension of the data set the number of cells would decrease to unity and the mean span would become one. A span mean of one is a lower limit on the y-axis. The hyperbolic nature follows the earlier argument that in the absence of boundary relocations during aggregation the distributions would change by simply dividing the resolution axis values by the resolution number. This would be equivalent to the product of resolution number and the mean span at that resolution equalling a constant (the span mean for resolution one).

Boundary relocations do occur during aggregations, however. As the resolution number increases and eventually surpasses the mean span of the reference data set, many smaller map units and even the thinner protrusions of the larger map units begin to disappear as they fail to dominate the larger cells. Elimination of small spans will immediately create large spans and the mean span can be expected to increase with increasing resolution number.

For the data sets of the intensive study, the product of resolution number and mean span is plotted versus resolution number as shown in Figure 18. This product is a mean span distance since the changing cell dimension is being accounted for by the resolution number. The relationship is clearly not a constant over the range of resolution numbers. An increase in mean span due to boundary relocations and deletions is more than offsetting the decrease in mean due to increasing cell size alone. It is interesting that the net increase so closely fits a linear regression model.

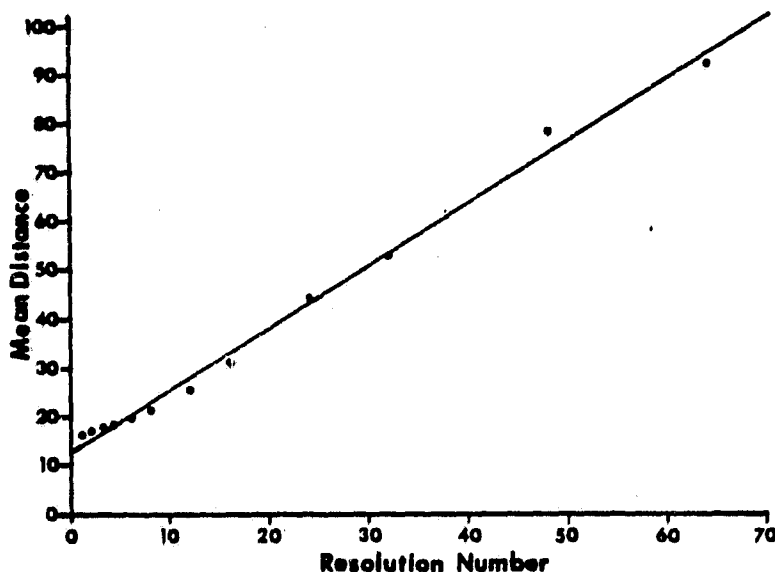


Figure 18. Mean Span Distance versus resolution number.

#### Data Characterization

A mathematical relationship for Figure 18 was not sought by standard statistical and/or curve fitting techniques since it was realized that such a relationship - even though clearly existent - would not likely hold for other map data sets. For example maps with multimodal span distributions would have mean distances that change in an altogether different pattern with changing resolution. This would quite possibly make the equivalent of Figure 18 non-linear and completely alter the relationship of errors to span distribution mean.

Since the span distribution characterizes the composition of map interboundary distances and spans correspond directly to the compact, sequential geocoding process, estimation of mapping error at



a certain cell size should be possible from the relative frequencies of the span distribution rather than the mean span. Various distributions including multimodal possibilities may arise from different map structures and yet have identical mean spans. The span distributions, however, will represent such differences in map structure.

## MAPPING ERROR AND THE SPAN DISTRIBUTION

### A POSITIONAL AVERAGE MODEL

The system performance experiments demonstrated that even though individual mapping unit error for varying sizes and shapes may behave erratically with changing cell size, an average effect over the map is quite well behaved. Earlier the mean span distance was found to logically relate to changing cell size. It was noted, however, that this relationship was probably unique to the particular data set being analyzed and the result is not applicable to other data sets.

The usefulness of the span distribution itself to overcome the shortcomings of the mean span was noted. The practical relevance of the span distribution to predicting mapping error depends on the definition of a relationship between them. Further discussion will introduce the pursuit of such a relationship.

### The Size and Orientation of Cells

Several conceptual relationships between the span distribution and mapping error at a given cell size can be visualized. All depend

on the span as a discrete interboundary distance and the cell size which must quantize that distance.

Consider one spatial dimension, an isolated span of  $n$  spatial units and a cell dimension of  $m$  spatial units. Let the cellularization begin in alignment with one end of the span. If  $m$  equals  $n$ , the cell would represent the span exactly and there would be no mapping error; if the  $m$  unit cell becomes greater than  $2n$ , there would be 100% mapping error; for  $m$  between 1 and  $2n$  the error depends on the divisibility of the span by the cell.

Further consider  $m=n$  with end alignment. There would be no mapping error. If the spatial relationship were not end aligned but rather left to chance, there could result up to 49% mapping error in the cases where two cells join near the midpoint of the span.

For a particular cell dimension of  $m$  units the factors which influence the error in mapping a span of length  $n$  units are the chance orientation of the cell or cells with respect to the span, and the  $m$ - $n$  size relationship. In addition the interdependency of spans in the two spatial dimensions and the interplay of boundary adjustments between adjacent sequential spans of a map transect may warrant investigation.

#### The Span Distribution and Mapping Error

A conclusion was that interboundary distance cannot be usefully summarized by any parameter (region area) or statistic (mean span) which is not uniquely descriptive of the interboundary distance distribution. The distribution itself is uniquely characteristic of the corresponding map data.

Recall that the span distribution is discrete. Only fixed integer multiples of some linear spatial unit are present. Actual interboundary distances are continuous and may take on any real value. The span distribution must be derived with a sufficiently small spatial unit so as to avoid quantization impact on the interboundary distance characteristics of the map. "Sufficiently small" is not impossible to evaluate nor is it so small as to be prohibitively impractical.

The span distribution has ordinates of relative frequency versus abscissas of discrete span size. A graph of mapping error versus changing cell size has ordinates of percent error and abscissas of cell size. Earlier discussions suggested that physical relationships of cell size and span size are the source of mapping error. What is needed is a matrix of error components for  $m$  values (cell sizes) versus  $n$  values (span sizes). This matrix multiplied by the span distribution vector would yield the error prediction vector. Figure 19 diagrams this proposed relationship between the span distribution and the mapping error. Mathematically

$$\bar{e}(m) = \bar{f}(n) \cdot \bar{g}(n,m)$$

The  $\bar{g}(n,m)$  matrix would represent universal physical relationships between various cell and span sizes in terms of a corresponding mapping error. The span distribution  $\bar{f}(n)$  would represent the unique character of the particular map---the actual combination of spans which comprise the map. The matrix product will yield a vector of error fractions versus cell sizes  $m$  for the map represented by  $\bar{f}(n)$ .

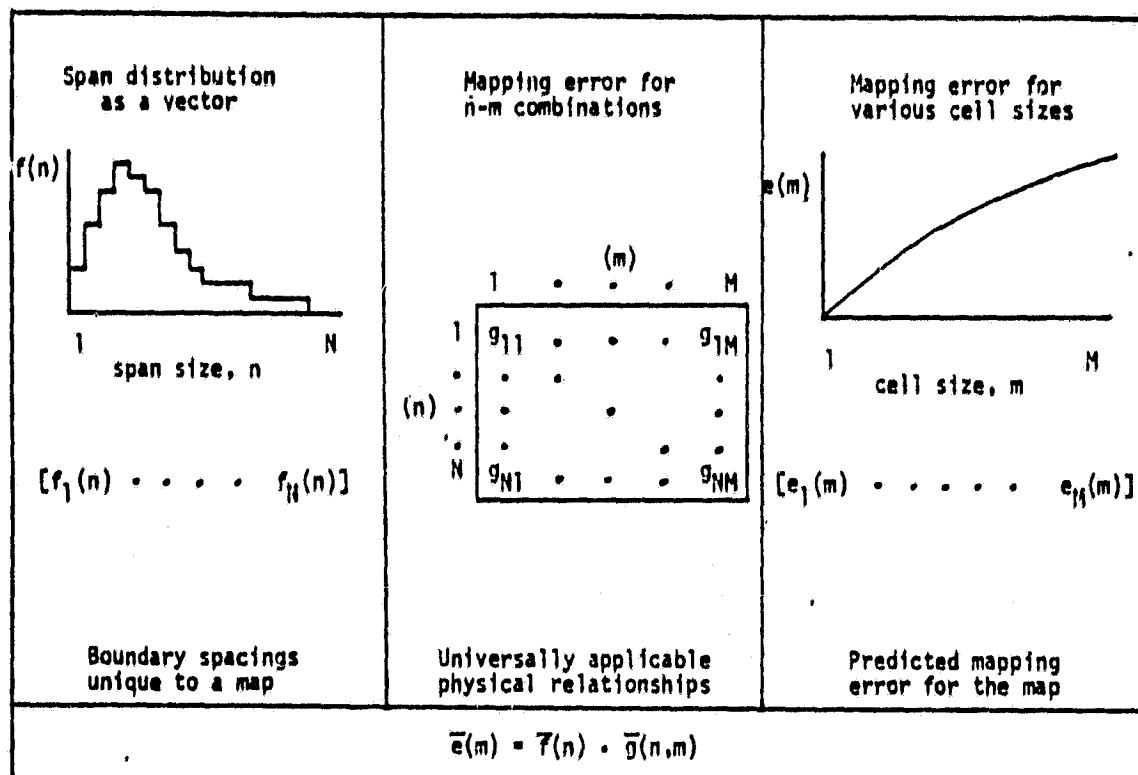


Figure 19. Components of the proposed mathematical relationship between span distribution and mapping error vectors.

### A Positional Average Model

The problem at hand is to determine  $\bar{g}(n,m)$ . Earlier experiments pointed to the chance orientation of cells with respect to spans as a source of wide variation in mapping error. The experiments were encouraging, however, in the discovery of a well behaved positional average error. For this reason a positional average error fraction was considered a reasonable starting model for  $\bar{g}(n,m)$ .

The model was pursued by manually diagramming all cell positions for each n-m combination (span size-cell size). Enough combinations were diagrammed to observe a pattern or function which could be

employed to generate the  $\bar{g}(n,m)$  matrix regardless of the dimensions  $N,M$ . Assumptions were a one-dimensional basis and a homogeneous background. Figure 20 presents a few examples of the procedure. Consider the first case  $n=2, m=2$ . There are two ( $m=2$ ) possible positions for the cells. Position one aligns exactly with the span and mapping error would be zero. Position two is more complex, however. There are two cells, each which might result in one half of the span being coded incorrectly or 50% error. By the cell dominance coding rule in this case with a 50-50 split of the cell there is a random choice made between the class of the span and the class of data outside the span. The calculation is (2 cells) ( $\frac{1}{2}$  chance of error)(50% error)---which is the 50% tabulated.

For the case  $n=2, m=4$ , position one, the entire span, 100%, might be error if the cell is coded to the data class outside the span. This can happen on the random tie breaker with a 50% chance. Hence  $\frac{1}{2}$  chance of 100% mapping error is the 50% tabulated. From the case  $n=2, m=5$  it should become obvious that all  $m > 2n$  are 100% average error since the span is never capable of dominating the cell.

The procedure sketched in Figure 20 was continued for  $n=1,2, \dots, 10$  and  $m=2,3, \dots, 6$ . Average error was tabulated as fractions in reduced form. Observation of the column  $m=2, n=1,2, \dots, 10$  revealed an obvious pattern in denominators which led to expression of all denominators as the value of the product  $m$  times  $n$ . Table 3 contains these fractions for  $n=1,2, \dots, 10$  and  $m=2,3, \dots, 9$ . The denominator pattern has been noted. Numerators are constant within a column except for the entries in parentheses. As mentioned many times

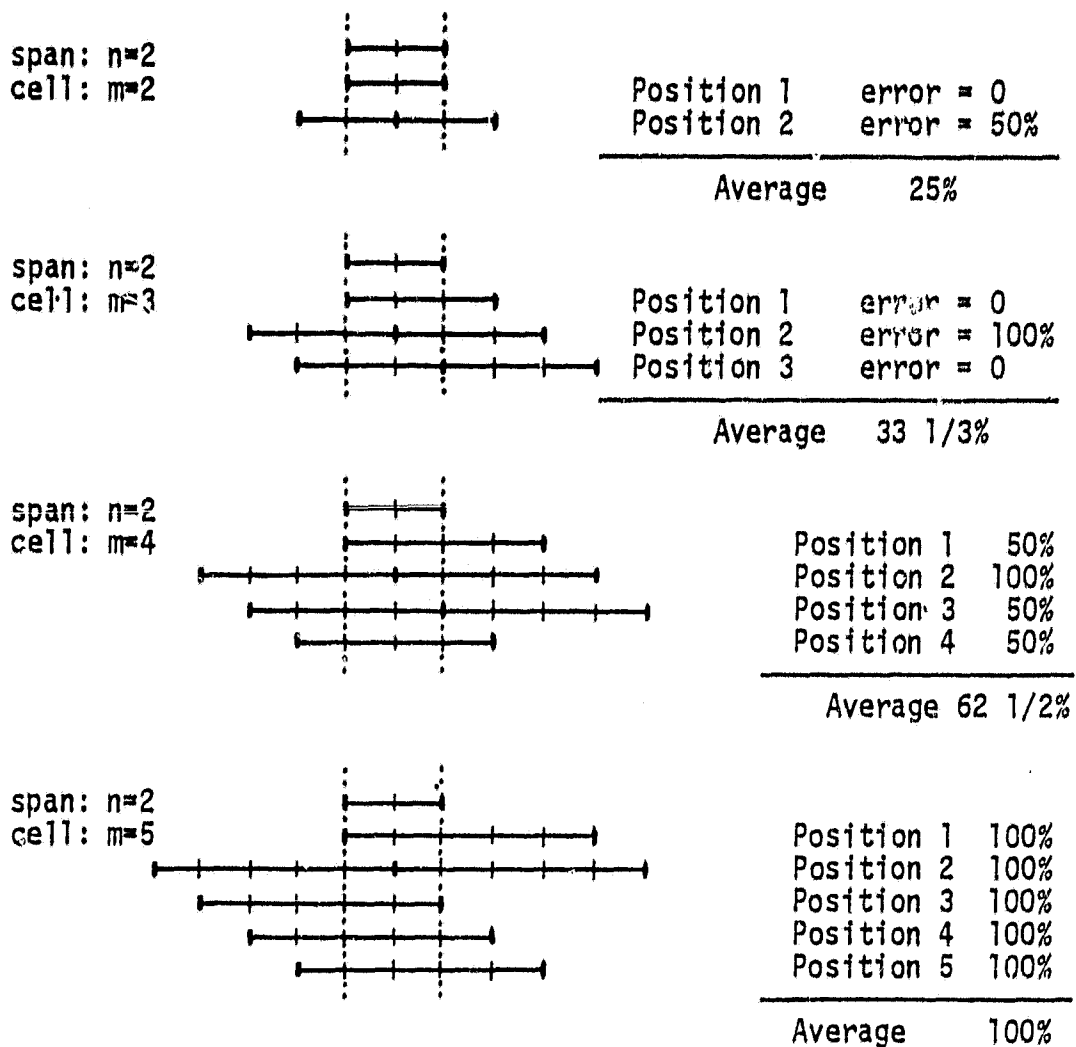


Figure 20. The procedure for observing  $\bar{g}(n,m)$  by systematic sketches of the  $m$  positions of a cell of size  $m$  with respect to a span of size  $n$ . See text for discussion.

Table 3. Positional Average Error Fractions as a Model of  $\bar{g}(n,m)$ .

$\begin{smallmatrix} m \\ n \end{smallmatrix}$	2	3	4	5	6	7	8	9
1	1/2	1	1	1	1	1	1	1
2	1/4	2/6	(5/8)	1	1	1	1	1
3	1/6	2/9	4/12	6/15	(12/18)	1	1	1
4	1/8	2/12	4/16	6/20	9/24	12/28	(22/32)	1
5	1/10	2/15	4/20	6/25	9/30	12/35	16/40	20/45
6	1/12	2/18	4/24	6/30	9/36	12/42	16/48	20/54
7	1/14	2/21	4/28	6/35	9/42	12/49	16/56	20/63
8	1/16	2/24	4/32	6/40	9/48	12/56	16/64	20/72
9	1/18	2/27	4/36	6/45	9/54	12/63	16/72	20/81
10	1/20	2/30	4/40	6/50	9/60	12/70	16/80	20/90

previously, when  $m$  becomes greater than twice  $n$ , the error becomes 100%.

The unity entries can be mathematically predicted, the remaining denominators mathematically predicted, a numerator pattern versus  $m$  can be mathematically described and even the small additive factor for the numerators of the same terms in parenthesis fit a pattern.

The entire matrix can be generated to any dimensions required by the following rules:

- (1) for  $m > 2n$ ,  $\bar{g}(n,m) = 1$
- (2) for  $m \leq 2n$  all denominators are  $mn$
- (3) for  $m \leq 2n$  numerators fit the pattern  $\text{num}(m) = \text{num}(m-1) + [m/2]$
- (4) additive numerator corrections for  $m = 2r$  are

$$\sum_{i=1}^{n-1} i$$

Some physical bases for the entries may be noted. The source of the unity entries is physically obvious. The denominator  $m$  times  $n$  arises from averaging over  $m$  positions and from the mapping error being part or all of the  $n$ -unit span expressed as a fraction of  $n$ . The additive factors in the terms where  $m=2n$  arise from random tie breaking in assignment of cell dominance. A physical reason for the remaining numerator sequence is not immediately apparent.

#### Prediction of Mapping Error

With the rules defined for generating  $\bar{g}(n,m)$  as a positional average error fraction matrix, it was a reasonably simple task to generate the matrix on a digital computer. The span distribution obtained from the base data set of the performance experiment encompassed spans to 120 units in length and the aggregations used in the study of changing cell size ranged from 2 to 64. Hence the  $\bar{g}(n,m)$  matrix generated was 120 elements long ( $n$ ) and 64 elements wide ( $m$ ). The span distribution,  $\bar{f}(n)$ , a 120 element vector, was then multiplied by  $\bar{g}(n,m)$  to predict the mapping error vector  $\bar{e}(m)$  in accordance with Figure 19.

The predicted mapping error and the experimental mapping error for cell sizes up through 64 spatial units are compared in Figure 21. There is a consistent over-estimation of mapping error by consideration of positional average relationships between spans and cells in a one dimensional model.



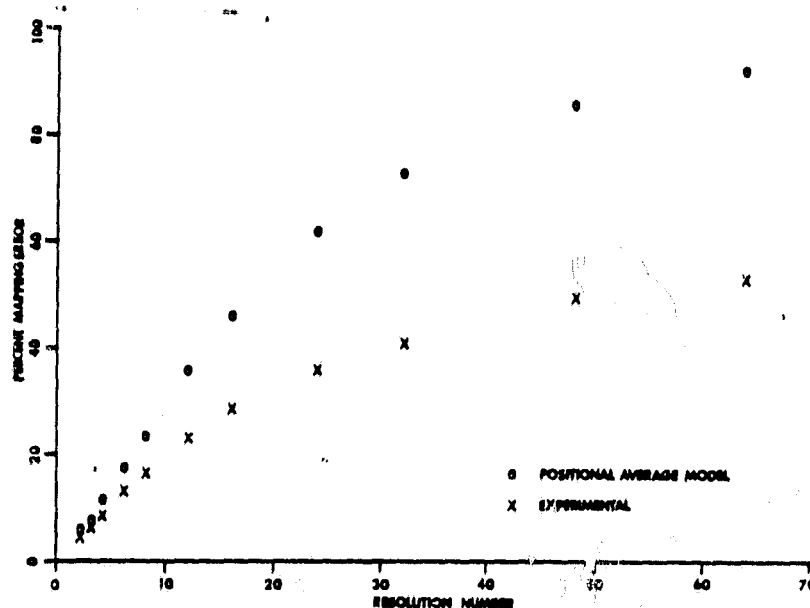


Figure 21. Predicted versus experimental mapping error.  $\bar{g}(n,m)$  used was the positional average error fraction matrix.

#### MAPPING ERROR AND THE SPAN DISTRIBUTION

##### - CORRECTION FOR SPAN ADJACENCIES

The positional average model was derived by observing isolated spans of  $n$  spatial units and the  $m$  positions of a cell of  $m$  spatial units. Isolation was introduced by consideration of the span with respect to a uniform, extensive background. In each case mapping error was considered in terms of the  $n$  units interior to the span and no consideration was given to correct mapping of the background. Cell dominance determined whether or not the units of the span were correctly or incorrectly mapped. In actual map applications, spans occur in a sequence along any map transect. This adjacency of spans may influence cell dominance decisions and alter the corresponding

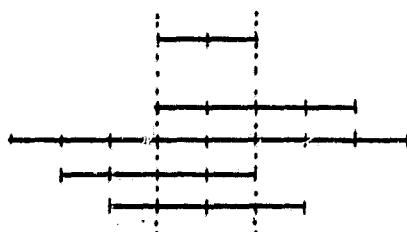
mapping error. An attempt to model such effects as a correction to the positional average mapping error matrix  $\bar{g}(n,m)$  was made.

### Span Adjacency Influence

Consider Figure 22. Part A demonstrates the calculation of positional average mapping error when the region surrounding the span is large and homogeneous. In part B an adjacent span of one unit is assumed to be different from the background data as well as different from the data associated with the two-unit span being analyzed. Only positions three and four from the calculation in part A are repeated to demonstrate the impact of adjacency. In position three the

span:  $n=2$

cells:  $m=4$



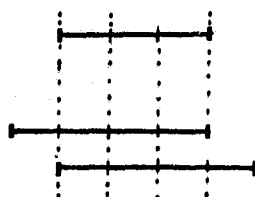
Position 1	50%
Position 2	100%
Position 3	50%
Position 4	50%
Average	62½%

A) the one dimension, positional average calculation.

spans:  $n=1$  and  $n=2$  adjacent

cell:  $m=4$

positions 3 & 4



Position 3	0%
Position 4	50%

B) the one dimensional, positional average calculation when a one unit span is known to be adjacent (A reconsideration of positions 3 and 4 of part A)

Figure 22. Impact of an adjacent span on calculation of positional average mapping error.

presence of the one-unit span causes the data class of the two-unit span to dominate the cell. The cell would be mapped as the data class of the two-unit span and no mapping error (in terms of the span itself) would occur. In part A the same cell position stood a 50-50 chance of totally misrepresenting the span--a 50% error. Clearly for position three, adjacency considerations could reduce the positional average mapping error. The example does demonstrate the desired effect -- a reduction in mapping error magnitude when adjacency effects are considered.

#### Describing Adjacent Span Corrections

Considering Figure 22 part B, position three in comparison to part A, position three, there is a 50% error term removed. This correction is only valid, however, for the case where a one-unit span occurs adjacent to the two-unit span being analyzed. At present, there is no evidence pointing to interdependency of spans. There is no cause-effect relationship known that would dictate the relative frequencies of various spans occurring adjacent to a known span. The assumption is made that the likelihood of an adjacent span of one unit is estimated simply by the relative frequency of one-unit spans in the map, i.e.  $f(1)$ . Simply stated the assumption is that spans are independent of one another. The 50% error correction for position three should then be weighted by the likelihood of the one-unit span actually occurring. The correction is then described as  $\frac{1}{2}f(1)$ .

The potential occurrence of multiple spans of corrective influence must also be considered. In Figure 23, the position of the two-unit

span with respect to the six-unit cell will result in 100% error if potential adjacent span influences are ignored.

Span: 2 units

Cell: 6 units

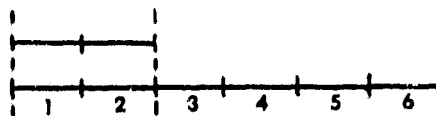


Figure 23. One possible position of a six-unit cell with respect to a two-unit span.

The two unit span will be correctly mapped if immediately following the two-unit span there are three single-unit spans,  $f^3(1)$ . If there followed two two-unit spans, a three way tie would be broken by random chance and a one-third chance of correct mapping would arise from the two two-unit spans,  $1/3f^2(2)$ . In both of these cases, multiple spans of different content are being used as the corrective factor. According to a previous assumption, an adjacent span is only different from the background data and the data associated with the span it is next to. In looking at the three single-unit spans, there is no way to insure that the contents of positions two and four (or positions two and five, three and five) will be different. The same is true for the case of the two two-unit spans. Therefore, any multiple spans will not be included in the derivation of correction effects. Also, using the same reasoning, any cross products of spans, such as  $f(1)f(2)$  or  $f(2)f(4)f(5)$ , will not be included. The only possible correction factor, then for figure 23 is a single two-unit span,  $f(2)$ , which will give a three-way tie and a resulting correction term of  $1/3f(2)$ .

Another situation that has to be considered is the occurrence of more than one way to correct the positional error as illustrated in Figure 24.

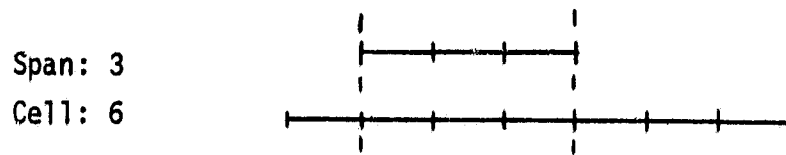


Figure 24. One possible position of a six-unit cell with respect to a three-unit span.

A one-unit span,  $f(1)$ , will correct the error 100%; a two-unit span,  $f(2)$ , will also correct the error 100%, as will a three-unit span,  $f(3)$ , a four-unit span,  $f(4)$ , etc. ... In general the series  $f(1) + f(2) + f(3) + \dots$  will correct the positional error for this case.

In the last example, the adjacent spans  $f(1)$ ,  $f(2)$ ,  $f(3)$ , ... could be placed either to the right or to the left of the three-unit span; the resulting correction was the same for each case. However, in the example illustrated by Figure 25, the placement of the adjacent spans does influence the resulting correction.

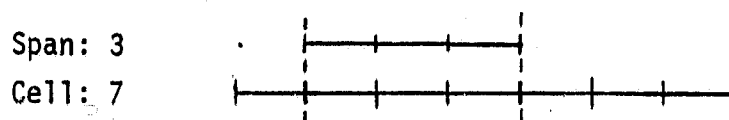


Figure 25. One possible position of a seven-unit cell with respect to a three-unit span.

A one-unit span on either the right or the left side will give the three unit span a 50-50 chance of being mapped correctly. A two-unit span on the right side result in correctly mapping the three-unit span, but a two-unit span on the left side results in a 50-50

chance of the three unit span being mapped correctly. At this point a decision had to be made concerning placement effects. In trying to keep the correction model from not becoming too complex, it was decided to assume that all possible placements of adjacent spans within the cell must occur before any series terms ( $f(n_1) + f(n_2) + f(n_3) + \dots$ ) can be considered. Also, whenever an adjacent span could occur on either side with the same results, its probability was treated as if it could only occur in one of the positions. Therefore, in this example,  $f(1)$  will give a 50-50 chance of being correct,  $f(2)$  will fully correct, and  $f(3) + f(4) + f(5) + \dots$  will give a 50-50 chance of being correct. The series  $f(1) + f(2) + f(3) + \dots$  on the left side would also give a 50-50 chance of being correct, but this is an incorrect expression since it does not follow the established assumptions. Correction term rules and notation is, therefore, established.

#### The Correction Matrix

There should exist a matrix  $\bar{h}(n,m)$  corresponding one-to-one with the derived  $\bar{g}(n,m)$ . The entries of  $\bar{h}(n,m)$  are the subtractive correction factors which account for reduced mapping error with adjacency influence considered. To determine this matrix, observations of spans of  $n$  units in combination with cells of  $m$  units were again made as was done in the positional average analysis reported earlier. Using the established correction notation, average correction terms must be developed for each entry of  $\bar{h}(n,m)$ .

This procedure was pursued with hope that a generating pattern could again be defined. Some typical correction expressions obtained are organized in Table 4. From the discussion of Figures 23 through 25 and the entries below the lower dotted line in Table 4, it should be apparent that the larger the cell size for a given span size, the greater the number of correction terms. This arises from the larger cell sizes having positions with respect to a span which extend well beyond the span. The further cell extension beyond the end of the span provides opportunity for series terms or larger single spans to have a corrective influence on mapping error. Also, the zeroes above the upper dotted line are the result of the cell-to-span size ratio becoming too large and, therefore, a span having no chance of being mapped correctly.

Table 5 shows a general pattern of the entries in Table 4; three distinct regions are evident. Below the lower dotted line is the "regular" region, denoted by  $x$ , between the dotted lines is the "transition" region denoted by  $T$  and  $K$ , and above the upper dotted line is a region of zeroes, denoted by  $0$ . The basic pattern of the regular region extends throughout the transition region, with the exception of an extra corrective term (in brackets) for each entry.

Further inspection of Table 4 reveals that  $n$  is simply a divisor of all terms in the overall pattern (excluding the extra corrective terms in the transition region). Therefore, the tabulation of these terms should actually be of  $m$ , the cell size, and various  $f(n)$  which contribute corrective effects; Table 6 represents this reorganization. The  $m$  times  $n$  in the denominator is understood and, therefore, dropped

Table 4. Correction expressions for selected m-n combinations.

m \ n	3	4	5	6	7	8
2	$\frac{2}{6}f(1) + \frac{1}{2}f(1) + \frac{1}{2}f(2)$	$\frac{2}{3}f(1) + \frac{2}{3}f(2) + \frac{1}{3}f(3)$	$\frac{2}{3}f(1) + \frac{1}{3}f(2) + \frac{1}{3}f(3)$	$\frac{11}{6}f(1) + \frac{5}{6}f(2) + \frac{1}{3}f(3)$	0	0
3	$\frac{2}{9}f(1) + \frac{2}{9}f(2)$	$\frac{2}{3}f(1) + \frac{2}{3}f(2)$	$\frac{2}{3}f(1) + \frac{2}{3}f(2)$	$\frac{13}{9}f(1) + \frac{13}{9}f(2) + \frac{1}{3}f(3)$	$\frac{2}{3}f(1) + \frac{5}{3}f(2) + \frac{2}{3}f(3)$	$\frac{2}{3}f(1) + \frac{2}{3}f(2) + \frac{2}{3}f(3) + \frac{1}{3}f(4)$
4	$\frac{2}{12}f(1) + \frac{2}{12}f(2)$	$\frac{2}{4}f(1) + \frac{2}{4}f(2)$	$\frac{2}{4}f(1) + \frac{2}{4}f(2)$	$\frac{13}{12}f(1) + \frac{13}{12}f(2)$	$\frac{3}{4}f(1) + \frac{5}{4}f(2) + \frac{3}{4}f(3)$	$\frac{1}{4}f(1) + \frac{1}{4}f(2) + \frac{1}{4}f(3) + \frac{1}{4}f(4)$
5	$\frac{2}{15}f(1) + \frac{2}{15}f(2)$	$\frac{2}{5}f(1) + \frac{2}{5}f(2)$	$\frac{2}{5}f(1) + \frac{2}{5}f(2)$	$\frac{13}{15}f(1) + \frac{13}{15}f(2)$	$\frac{2}{5}f(1) + \frac{5}{5}f(2) + \frac{2}{5}f(3)$	$\frac{2}{5}f(1) + \frac{2}{5}f(2) + \frac{2}{5}f(3)$
6	$\frac{2}{18}f(1) + \frac{2}{18}f(2)$	$\frac{2}{6}f(1) + \frac{2}{6}f(2)$	$\frac{2}{6}f(1) + \frac{2}{6}f(2)$	$\frac{13}{18}f(1) + \frac{13}{18}f(2)$	$\frac{2}{6}f(1) + \frac{5}{6}f(2) + \frac{2}{6}f(3)$	$\frac{2}{6}f(1) + \frac{2}{6}f(2) + \frac{2}{6}f(3)$
7	$\frac{2}{21}f(1) + \frac{2}{21}f(2)$	$\frac{2}{7}f(1) + \frac{2}{7}f(2)$	$\frac{2}{7}f(1) + \frac{2}{7}f(2)$	$\frac{13}{21}f(1) + \frac{13}{21}f(2)$	$\frac{2}{7}f(1) + \frac{5}{7}f(2) + \frac{2}{7}f(3)$	$\frac{2}{7}f(1) + \frac{2}{7}f(2) + \frac{2}{7}f(3)$
8	$\frac{2}{24}f(1) + \frac{2}{24}f(2)$	$\frac{2}{8}f(1) + \frac{2}{8}f(2)$	$\frac{2}{8}f(1) + \frac{2}{8}f(2)$	$\frac{13}{24}f(1) + \frac{13}{24}f(2)$	$\frac{2}{8}f(1) + \frac{5}{8}f(2) + \frac{2}{8}f(3)$	$\frac{2}{8}f(1) + \frac{2}{8}f(2) + \frac{2}{8}f(3)$



Table 5. Pattern of correction expressions of Table 4.

n	m																		
		3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18		
2		X	K	K	T	0	0	0	0	0	0	0	0	0	0	0	0		
3		X	X	X	K	K	T	T	0	0	0	0	0	0	0	0	0		
4		X	X	X	X	X	K	K	T	T	0	0	0	0	0	0	0		
5		X	X	X	X	X	X	X	K	K	T	T	T	0	0	0	0		
6		X	X	X	X	X	X	X	X	X	K	K	T	T	T	T	T		
7		X	X	X	X	X	X	X	X	X	X	X	K	K	T	T	T		
8		X	X	X	X	X	X	X	X	X	X	X	X	X	K	K	T		
9		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	K		

Table 6. Coefficients of  $\bar{f}(n)$  for various  $m$  to generate first-order correction expressions for span adjacency effects.

$m$	$k$	$f(1)$	$f(2)$	$f(3)$	$f(4)$	$f(5)$	$f(6)$	$f(7)$	$f(8)$
3	1	$k - \frac{1}{3}$							
4	2	$k$							
5	2	$k$	$k$						
6	3	$k$	$(2k-1) - \frac{2}{3}$						
7	3	$k$	$2k$	$k$					
8	4	$k$	$2k-1$	$2k-1$					
9	4	$k$	$2k$	$(3k-1) - \frac{3}{3}$	$k$				
10	5	$k$	$2k-1$	$3k-2$	$2k-1$				
11	5	$k$	$2k$	$3k-1$	$3k-1$	$k$			
12	6	$k$	$2k-1$	$3k-2$	$(4k-4) - \frac{4}{3}$	$2k-1$			
13	6	$k$	$2k$	$3k-1$	$4k-2$	$3k-1$	$k$		
14	7	$k$	$2k-1$	$3k-2$	$4k-4$	$4k-4$	$2k-1$		
15	7	$k$	$2k$	$3k-1$	$4k-2$	$(5k-4) - \frac{5}{3}$	$3k-1$	$k$	
16	8	$k$	$2k-1$	$3k-2$	$4k-4$	$4k-4$	$4k-4$	$2k-1$	
17	8	$k$	$2k$	$3k-1$	$4k-2$	$5k-4$	$5k-4$	$3k-1$	$k$

(all entries require  $1/mn$ )

from the entries. The number pattern is expressed as functions of  $k$  where  $k$  is the largest integer in  $m/2$ , i.e.  $k = [m/2]$ .

Three distinct patterns emerge for the extra corrective terms in the transition region, these being the even and odd  $k$  terms, (with respect to  $m$ ), and the  $T$  terms (all from Table 5). Table 7 shows the even and odd  $k$  terms, and Table 8 shows the  $T$  terms.

Table 7. Terms for even and odd  $m$  in the "K" entries of Table 5.

K TERMS			
EVEN $m$		ODD $m$	
$m$	$K$	$m$	$K$
4	$1/2$	5	$-0/2f(1)+1$
6	$2/2$	7	$-1/2f(1)+f(2)+2$
8	$3/2$	9	$-2/2f(1)+f(2)+f(3)+3$
10	$4/2$	11	$-3/2f(1)+f(2)+f(3)+f(4)+4$
12	$5/2$	13	$-4/2f(1)+f(2)+f(3)+f(4)+f(5)+5$

The pattern in Table 7 is obvious, but the one in Table 8 is not, except for the pattern of the constants of each entry. After little success in finding an overall pattern for Table 8, it was decided to approximate the  $T$  terms by only including the constants.

Now that a useable pattern has been determined for all the entries of Table 4, the subtractive corrective matrix,  $\bar{h}(m,n)$ , will be developed. The entries of Table 6 are organized as a matrix which matches in sizes to the factors of  $\bar{f}(n)$  required to generate a correction expression; call this matrix  $\bar{c}(m,n)$ . Multiplying  $\bar{c}(m,n)$

Table 8. The "T" entries of Table 5.

m \ n	6	8	9	10	11	12
	$-\frac{11}{6}f(1)-\frac{5}{6}f(2)+\frac{1}{3}$					
2						
3		$-\frac{7}{3}f(1)-\frac{1}{3}f(2)$ $-\frac{1}{3}f(3)+1$	$-\frac{5}{3}f(1)-\frac{9}{3}f(2)$ $-\frac{4}{3}f(3)-\frac{3}{4}f(4)+\frac{1}{3}$			
4				$-\frac{13}{4}f(1)-\frac{3}{4}f(2)$ $-\frac{7}{4}f(3)-\frac{1}{4}f(4)+2$	$-\frac{9}{4}f(1)-\frac{14}{9}f(2)$ $-\frac{1}{2}f(3)-\frac{1}{2}f(4)$ $-\frac{5}{4}f(5)+1$	$-\frac{22}{12}f(1)-\frac{37}{12}f(2)$ $-\frac{52}{12}f(3)-\frac{24}{12}f(4)$ $-\frac{11}{4}f(5)+\frac{1}{3}$

by the vector span distribution  $\bar{f}(n)$  will yield a vector  $\bar{v}(m)$  which has single term expressions containing  $n$ . If  $\bar{v}(m)$  is expanded to two dimensions by inserting values of  $n$  into the rows, and if the corrective terms in the diagonal transition region are then added and the zero entries properly placed, then the result will be the desired  $\bar{h}(n,m)$ .

### Implementing the Correction

To implement the procedure outlined in the previous paragraph, some generating pattern must be found to create the  $\bar{c}(m,n)$  matrix equivalent to Table 6. The small adjustment fractions  $1/3$ ,  $2/3$ ,  $3/3$  etc. can be easily generated as  $n/3$  for entries at  $m=3n$  so they can be ignored temporarily in seeking a pattern. The form of all entries then becomes a coefficient of  $k$  and a subtractive number.

Table 9 summarizes the coefficients of  $k$  in the  $\bar{c}(m,n)$  matrix. This pattern is easy to define. For a column of  $\bar{f}(n)$  run zeros from  $m=1$  to  $m=2n$ , count upward from 1 at  $m=2n+1$  to  $m=3n$  and fill the column from  $m=3n$  to whatever limit is desired with the value at  $m=3n$ .

The generation of the subtractive numbers is not so obvious. Table 10 summarizes the subtractive numbers for the entries of  $\bar{c}(m,n)$  in two groups, one for  $m$  even and one for  $m$  odd. A pattern for generating entries in one group can be used for the other and the groups combined to result in the subtractive number table.

First consider the constant column entries below the lower dividing line at  $m=3n$ . The differences between columns fits the pattern 1,1,2,2,3,3 etc. The last line, largest value of  $m$  desired,

can be generated based on this difference pattern and the columns propagated upward to  $m=3n$ .

Table 9. Coefficients of K in the entries to the  $\bar{c}(m,n)$  matrix.

m	K	f(1)	f(2)	f(3)	f(4)	f(5)	f(6)	f(7)
3	1	1						
4	2	1						
5	2	1	.1					
6	3	1	2					
7	3	1	2	1				
8	4	1	2	2				
9	4	1	2	3	1			
10	5	1	2	3	2			
11	5	1	2	3	3	1		
12	6	1	2	3	4	2		
13	6	1	2	3	4	3	1	
14	7	1	2	3	4	4	2	
15	7	1	2	3	4	5	3	1
16	8	1	2	3	4	5	4	2
17	8	1	2	3	4	5	5	3

The wedge of entries between the dividing lines is generated from the difference sequence in entries down a column. The first entry is one in each column. This occurs for  $m=2$   $1 + 2$  if  $m$  is even or  $m=2$   $1 + 3$  if  $m$  is odd; the column heading is  $f(1)$ . The differences down the column are 3,5,7,9,etc.

Utilizing the observed patterns, an algorithm was written to generate the matrices equivalent to Tables 9 and 10. Multiplying the entries of Table 9 by  $k$  and subtracting the entries of Table 10 results in a numerical equivalent to Table 6 which was called  $\bar{c}(m,n)$ .

Table 10. Subtractive numbers for entries to the  $\bar{c}(m,n)$  matrix.

m	k	f(2)	f(3)	f(4)	f(5)	f(6)	f(7)	f(8)	f(9)
7	3	0							
9	4	0	1						
11	5	0	1	1					
13	6	0	1	2	1				
15	7	0	1	2	4	1			
17	8	0	1	2	4	4	1		
19	9	0	1	2	4	6	4	1	
21	10	0	1	2	4	6	9	4	1
23	11	0	1	2	4	6	9	9	4
25	12	0	1	2	4	6	9	12	9
27	13	0	1	2	4	6	9	12	16
4	2								
6	3	1							
8	4	1	1						
10	5	1	2	1					
12	6	1	2	4	1				
14	7	1	2	4	4	1			
16	8	1	2	4	6	4	1		
18	9	1	2	4	6	9	4	1	
20	10	1	2	4	6	9	9	4	1
22	11	1	2	4	6	9	12	9	4
24	12	1	2	4	6	9	12	16	9
26	13	1	2	4	6	9	12	16	16
28	14	1	2	4	6	9	12	16	20

This  $\bar{c}(m,n)$  was then multiplied by  $\bar{f}(n)$ , the span distribution, to obtain  $\bar{v}(m)$ . Since the  $\bar{v}(m)$  vector is developed from Table 6 entries which include the common factor  $1/mn$ , the vector can be expanded to two

dimensions by including values of  $n$  in the entries and using  $n$  as the second dimension. The resultant is  $\bar{h}(n,m)$  without the extra corrective terms in the transistor region and the zero entries. The addition of these entries is easily accomplished and the result is  $\bar{h}(n,m)$ , the desired correction matrix for  $\bar{g}(n,m)$ . The entries are all fractions or zeroes representing the over-estimate of mapping error in  $\bar{g}(n,m)$ .

#### • Applying the Correction

The correction may be handled two ways. The  $\bar{h}(n,m)$  may be subtracted from  $\bar{g}(n,m)$  and  $\bar{e}(m)$  generated. Alternatively  $\bar{f}(n)$  multiplied by  $\bar{h}(n,m)$  will generate an error correction  $\bar{e}_c(m)$ . This error magnitude can be subtracted from the error model vector  $\bar{e}(m)$ . The second approach was taken to enable observation of the span adjacency error correction terms  $\bar{e}_c(m)$ . Percentage mapping error for the positional average model, the span adjacency correction, the corrected model and the experiment itself are compared in Table 9. In Figure 26 the positional average error model and the experimental results are repeated from Figure 21 and the span-adjacency-corrected model included for comparison.

The span adjacency correction to the model appears to account for a portion of the deviation between the experimental results and the uncorrected model predictions. The corrected model also "parallels" the experimental results and the remaining derivation may be a result of the one-dimensional simplifications in the models thus far.



Table 11. Mapping error comparisons for the positional average model, the span adjacency correction, the corrected model and the experimental results.

cell size $m$	$\bar{e}(m)$ positional average	$\bar{e}_c(m)$ span adjacency correction	$\bar{e}(m) - \bar{e}_c(m)$ corrected model	observed experimental results
2	5.33	0	5.3	3.8
3	7.14	0	7.1	5.8
4	10.94	0.28	10.7	8.8
6	17.00	0.78	16.2	12.8
8	23.06	1.61	21.5	16.0
12	35.16	4.12	31.0	22.6
16	45.41	5.82	39.6	28.0
24	61.18	8.80	52.4	35.4
32	72.26	9.68	62.6	40.4
48	85.16	8.76	76.4	48.0
64	91.52	6.19	85.3	52.3

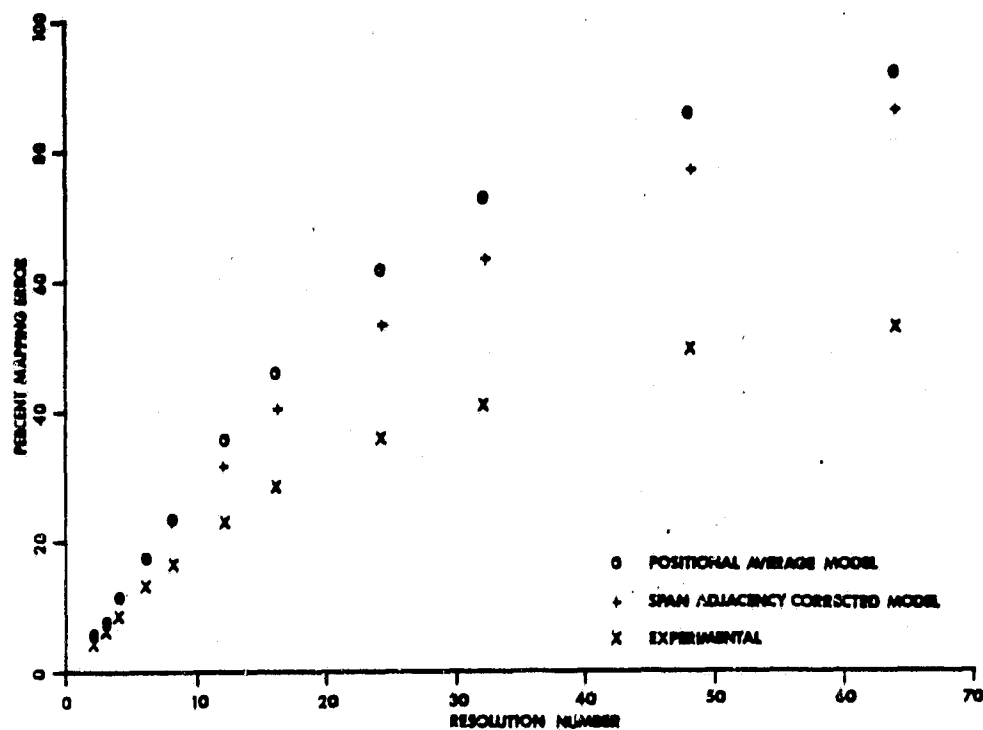


Figure 26. The experimentally observed mapping error compared to the positional-average model and the span-adjacency-corrected model.

Recall also that only the first order terms of span adjacency correction were utilized. If higher order terms and cross products had been included, the only possible effect would be an additional decrease in predicted mapping error ordinates throughout the cell size range. If the inclusion of higher order correction terms would bring the predicted error and experimentally observed error closer together in the cell size range 12 to 32, then the predicted error would probably far underestimate the experimentally observed error at larger cell sizes (beyond 32). This may be acceptable since two-dimensional interactions have not been included in the model.

## SUMMARY, CONCLUSIONS AND RECOMMENDATIONS

### Summary

The background provided on geoinformation systems outlined motivations, designs, comparisons, performance measurement techniques and suggested cell size selection bases for users of a grid system.

The basic and auxiliary data sets generated allowed mapping and inventory accuracy analyses as measures of system performance with changing cell size. An analysis of map characteristics was made. Two models of the relationship between a particular map and the system performance for that map were discussed.

### Conclusions

The following observations and conclusions are made as a result of the research reported in this paper:

(1) Although in the case of a single mapping unit, the behavior of mapping and inventory errors with increasing cell size may behave very erratically, an average affect over many mapping units in a map segment is apparently a well behaved non-decreasing function.

(2) The interboundary distance distribution does uniquely characterize a map.

(3) No summary parameter such as mean or mode of the distribution will ever provide a generally applicable prediction of performance since such parameters are not unique to a map.

(4) It appears that in the general case, the scanning of map distances in the two orthogonal directions corresponding to the cell dimensions, will yield distributions statistically dissimilar, hence implying that map sampling must be designed to randomly include transects in each dimension.

(5) A universal matrix model of the one-dimensional, average position of cells on a map does exist and may be generated to any desired dimensions as required.

(6) The influence of the adjacency of interboundary distances in one dimensional transects of a map is a significant, model-corrective factor. The first-order approximating model is also universal and generateable to whatever matrix dimensions are required.

(7) The total one-dimensional model of positional average corrected for first-order adjacency does not encompass enough of the physical interrelationships of cell size to map interboundary distances to predict mapping performance accurately.

(8) The existence of universally applicable model components in the studies undertaken is most encouraging to the continued development of the model matrix required to transform the inter-boundary distance distribution into a prediction of mapping accuracy.

#### Recommendations

The long range objective of the ongoing research in applications of a cellular information system is to define a procedure for selecting and justifying cell size. The composition of such a procedure has begun to crystallize. Sampling of map interboundary distances will estimate the distribution. A matrix model of physical relationships between cell sizes and map distances will transform that distribution into an estimate of mapping accuracy with various cell sizes. With knowledge of the operating costs of a particular processing system, AREAS or other cellular systems, the trade off between accuracy and cost for various cell sizes should enable cell selection to be made according to user needs.

The research reported and conclusions drawn are contributory to the understanding of underlying physical relationships which must exist for the outlined procedure to ever become reality. The gap between the accomplishments outlined and the selection procedure envisioned will be filled when the following recommendations have been successfully investigated:

(1) The modelling work should pursue two dimensional interactions in a fashion parallel to the sequence of investigations in the one dimensional case reported. This might require consideration of a joint or two-dimensional interboundary distance distribution. Higher order adjacency correction terms may also be required to bring predicted mapping error into closer agreement with experimental results. These pursuits would promote the derivation of the universal matrix model required.

(2) As a follow-up to model development as suggested in recommendation one, an analogous development could be pursued for inventory performance. Alternatively, inventory performance may be predictable from mapping performance when that model is finalized.

(3) A sampling technique, or at least guidelines, must be defined to enable a user to estimate the required form of interboundary distance distribution to drive the matrix model. No single distribution model will apply and, furthermore, the need for an estimate of the distribution (as opposed to an estimate of a parameter of a particular distribution) does not lend itself to known sampling techniques with defineable confidence intervals. For these reasons, this component of the desired overall selection procedure may take the form of empirical guidelines.

(4) Users of cellular systems, AREAS included, need to analyze operating costs in a manner that enables processing cost estimates for a given cell size. This is not always a trivial task inasmuch as some processes depend on several other parameters as well.

(5) Finally, the summary recommendation is made that the research be continued toward achievement of the cell size selection procedure. The work performed thus far indicates that such an approach is feasible and the requisite physical relationships do exist.

## BIBLIOGRAPHY

- [1] Joe D. Nichols, "Characteristics of Computerized Soil Maps" in Soil Science Society of America Proceedings, Vol. 39, pp. 927-932, 1975.
- [2] Phillip A. McDonald and Jerry D. Lent, "MAPIT-A Computer Based Data Storage, Retrieval and Update System for the Wildland Manager" in Proceedings of the 38th Annual Meeting of American Society of Photogrammetry, Washington, D.C., March 12-17, pp. 370-397, 1972.
- [3] W.A. Radlinski, "Modern Land Data Systems - A National Objective", Opening address to 1977 Annual Convention of American Society of Photogrammetry and American Congress of Surveying and Mapping, published - Photogrammetric Engineering and Remote Sensing, Vol. XLIII No. 7, pp. 887-890, July 1977.
- [4] "IRIS - Illinois Resources Information System", University of Illinois Center for Advanced Computations Feasibility Study - Final Report, Urbana, Illinois, 1972.
- [5] R.F. Tomlinson, "Geo-Information Systems and the Use of Computers in Handling Land Use Information" in Conference on Land Use Information and Classification, sponsored by Department of Interior of U.S. Geological Survey and the National Aeronautics and Space Administration, Washington, D.C., June 28-30, 1971.
- [6] Nevin A. Bryant and Albert L. Zobrist, "IBIS: A Geographic Information System Based on Digital Image Processing and Image Raster Datatype" in Symposium Proceedings of Machine Processing of Remotely Sensed Data, Purdue University, pp. 1A-1, 1A-7, June 29-July 1, 1976.
- [7] D. Steiner and T. Stanhope, "Data Base Development" Chapter 1 in Geographical Data Handling prepared for UNESCO/IGU Second Symposium on Geographical Information Systems by International Geographical Union Commission on Geographical Data Sensing and Processing, pp. 36-103, August, 1972.
- [8] R.F. Tomlinson, editor. Geographical Data Handling published for UNESCO/IGU Second Symposium on Geographical Information Systems by International Geographical Union Commission on Geographical Data Sensing and Processing, Ottawa, Ontario Canada, 1972.
- [9] Richard L. Phillips, "Computer Graphics in Urban and Environmental Systems", Proceedings of IEEE, Vol. 62, No. 4., pp. 437-452, April 1974.



- [10] Nevin A. Bryant and Albert L. Zobrist, "Integration of Socioeconomic Data and Remotely Sensed Imagery for Land Use Applications" in Proceedings of 2nd Annual Pecora Symposium sponsored by American Society of Photogrammetry and the U.S. Geological Survey, pp. 120-130, Oct. 25-20, 1976.
- [11] "A Land Classification Method for Land Use Planning" by the Land Use Analysis Laboratory at Iowa State University, 1973.
- [12] George Smith, Kris Van Gorkom, A.A. Dyer et al., Colorado Environmental Data Systems a final Report to the Colorado Department of Natural Resources by the College of Forestry and Natural Resources, Colorado State University, Fort Collins, Colorado, 1973.
- [13] D. Sinton. "Introduction to Spatial Data Manipulation and Analysis" Chapter 8 in Geographical Data Handling prepared for UNESCO/IGU Second Symposium on Geographical Information Systems by International Geographical Union Commission on Geographical Data Sensing and Processing, p. 719, Aug. 1972.
- [14] P. Switzer. "Estimation of the Accuracy of Qualitative Maps" in Display and Analysis of Spatial Data, NATO Advanced Study Institute, edited by John C. Davis and Michael McCullagh, John Wiley and Sons, New York, 1975.
- [15] Michael R. Hord and William Brooner "Land-Use Map Accuracy Criteria" in Photogrammetric Engineering and Remote Sensing, Vol. 42, No. 5, pp. 671-677, May 1976.
- [16] Charles R. Meyers Jr., Richard C. Durfee and Thomas Tucker, "Computer Augmentation of Soil Survey Interpretation for Regional Planning Applications" Oak Ridge National Laboratory Report ORNL-NSF-EP-67, April 1974.
- [17] Joe D. Nichols and Lindo J. Bartelli, "Computer-Generated Interpretive Soil Maps" in Journal of Soil and Water Conservation, Vol. 29(5), pp. 232-235, 1974.